

Nonparametric estimation of dynamics of monotone trajectories

Debashis Paul, Jie Peng and Prabir Burman¹

Department of Statistics, University of California, Davis

Abstract

We study a class of nonlinear nonparametric inverse problems. Specifically, we propose a nonparametric estimator of the dynamics of a monotonically increasing trajectory defined on a finite time interval. Under suitable regularity conditions, we prove consistency of the proposed estimator and show that in terms of L^2 -loss, the optimal rate of convergence for the proposed estimator is the same as that for the estimation of the derivative of a trajectory. This is a new contribution to the area of nonlinear nonparametric inverse problems. We conduct a simulation study to examine the finite sample behavior of the proposed estimator and apply it to the Berkeley growth data.

Keywords: autonomous differential equation; nonlinear inverse problem; monotone trajectory; nonparametric estimation; perturbation theory; spline

1 Introduction

Monotone trajectories describing the evolution of state(s) over time appear widely in scientific studies. The most widely studied are probably growth of organisms such as humans or plants (Milani, 2000; Erickson, 1976; Silk and Erickson, 1979). There are many parametric models for describing the features of growth curves, particularly in human growth (Hauspie et al., 1980; Milani, 2000). Most of these works focus on modeling the trajectories themselves or modeling the rate of change, i.e., the derivative of the trajectories. Other examples of monotone trajectories appear in population dynamics under negligible resource constraints (Turchin, 2003), in dose-response analysis in pharmacokinetics (Kelly and Rice, 1990), in auction price dynamics in e-Commerce (Jank and Shmueli, 2006; Wang et al., 2008; Liu and Müller, 2009), and in analysis of trajectories of aircrafts after take-off (Nicol, 2013). Some of these works are looking at function estimation with monotonic constraints and some of them are taking a functional data analysis approach.

In contrast, our goal here is to estimate the functional relationship between the rate of change and the state, i.e., the dynamics of the trajectory, through a nonparametric model. Many systems such as growth of organisms or economic activity of a country/region are intrinsically dynamic in nature (cf. Ljung and Glad, 1994). A dynamics model provides a mechanistic description of the system rather than a purely phenomenological one. Moreover, due to insufficient scientific knowledge, quite often there is a need for nonparametric modeling of the dynamical system. In addition, nonparametric fits can be used to develop measures of goodness-of-fit for hypothesized parametric models.

There is a large literature in modeling continuous time smooth dynamical systems through systems of parametric differential equations (see, e.g., Perthame, 2007, Strogatz, 2001). These methods have been used to model HIV dynamics (Wu, Ding and DeGruttola, 1998; Wu and Ding, 1999; Xia, 2003; Chen and Wu, 2008a, 2008b), the dynamic behavior of gene regulation networks (Gardner *et al.*, 2003; Cao and Zhao, 2008), etc. Approaches for fitting a parametric dynamics model include the maximum likelihood or nonlinear least squares. A recent approach proposed by Ramsay *et al.* (2007) and Cao *et al.* (2008) for parametric ordinary differential equations is based on the idea of balancing the model fit and the goodness of fit of the trajectories simultaneously.

¹Paul's research is partially supported by the NSF grants DMR-10-35468 and DMS-11-06690. Peng's research is partially supported by the NSF grant DMS-10-01256. Burman's research is partially supported by the NSF grant DMS-09-07622.

Another popular approach to fit dynamics models is a two-stage procedure (Chen and Wu, 2008a, 2008b; Varah, 1982), where the trajectories and their derivatives are first estimated nonparametrically and then the dynamics is fitted by regressing the fitted derivatives to the fitted trajectories. The two-stage approach can be easily adapted to estimate a nonparametric dynamics model. However, their performance is unsatisfactory due to the difficulty of resolving the bias-variance trade-off in a data dependent way. Brunel (2008) gives a comprehensive theoretical analysis of such an approach. Very recently, Hall and Ma (2014) proposed a one-step estimation procedure that mitigates some of the inefficiencies of two-stage estimators. However, this approach does not seem to extend naturally to estimate nonparametric dynamical systems.

There is also an extensive literature on the nonparametric estimation of monotone functions, e.g., Brunk (1970), Wright and Wegman (1980), Mammen (1991), Ramsay (1988, 1998). However, most methods in this field are not concerned with the estimation of the gradient function, except for Ramsay (1998) where the unknown function is modeled in terms of a second order differential equation and a smoothed estimate of its gradient is obtained as a byproduct.

A key observation of estimating the dynamics of monotone trajectories is that for any smooth monotone trajectory, its dynamics can be described by a first order autonomous differential equation. Specifically, if $X(t)$ is positive, strictly monotone and differentiable on a finite time interval, then we can express

$$X'(t) = (X' \circ X^{-1})(X(t)) = g(X(t)), \quad t \in [0, 1] \quad (1)$$

where $g = X' \circ X^{-1}$ is the gradient function. In this paper, we estimate the unknown gradient function g nonparametrically from discrete noisy observations of X . Specifically, we model the gradient function by a basis representation where the number of basis functions grow with the sample size. We adopt a nonlinear least squares framework for model fitting. We then carry out a detailed theoretical analysis and derive the rate of convergence of the proposed estimator.

We now highlight the major contributions of this work. Although there is a large literature on linear nonparametric inverse problems (Cavalier *et al.*, 2004; Cavalier, 2008; Donoho, 1995; Johnstone *et al.*, 2004), especially on the nonparametric estimation of the derivative of a curve (Gasser and Müller, 1984; Müller *et al.*, 1987; Fan and Gijbels, 1996), there is little theoretical development on nonlinear nonparametric inverse problems. Thus, our work makes a new contribution to this important area. In this paper, we first quantify the degree of ill-posedness of the estimation of the gradient function g as the number of basis functions grow to infinity. We then use this result to show that if g is p times differentiable then the L^2 -risk of the proposed estimator has the same optimal rate of convergence, viz., $O(n^{-2p/(2p+3)})$, as that of the estimator of the derivative of a trajectory assuming that the latter is $p + 1$ times differentiable. In Section 7, we show that the optimal rate of the proposed estimator is indeed the minimax rate for estimation of g under L^2 loss if the class of estimators is restricted to be uniformly Lipschitz. In the rest of the paper, unless otherwise specified, the phrase “optimal rate” refers to the best rate of convergence of the proposed estimator.

Among the few instances of nonparametric modeling of the gradient function known to us, Xue, Miao and Wu (2010) dealt with a related but different problem of estimating a parametric ODE with time-varying parameters, where the latter are modeled as unknown smooth functions of time. In a work most closely related to ours, Wu *et al.* (2014) proposed a sparse additive model for describing the dynamics of a multi-variate state vector and developed a combination of two-stage smoothing and sparse penalization for fitting the model. Their model can be seen as a multi-dimensional generalization of the autonomous ODE model studied here. In their paper, while deriving the risk bounds, it is assumed that whenever the gradient function g is p times differentiable, the state X is at least $3p + 1$ times differentiable. However, due to the representation $g = X' \circ X^{-1}$, it follows that g is p times differentiable if and only if X is $p + 1$ times differentiable. Therefore, at least for the one-dimensional state variable case, the assumptions made in Wu *et al.* (2014) are not satisfied in reality if p indeed denotes the maximal order of smoothness of g . This indicates that the rate

of convergence their estimator of g is not optimal for the current problem. It is also instructive to note that, due to the assumption about the additional degree of smoothness of the state variable, Wu et al. (2014) did not encounter the technical challenge posed by the ill-posedness of the problem.

The rest of the paper is organized as follows. In Section 2, we briefly describe the model and the estimation procedure. We present the main theoretical results in Section 3 and outline the main steps of the proof in Section 4. We present a simulation study in Section 5 and an application to the Berkeley growth data in Section 6. We discuss the optimality of the estimation of g in Section 7. Some proof details are provided in the Appendix (Section 8). Some derivations and graphical summaries are provided in the Supplementary Material (SM).

2 Model

The class of models studied in this paper is of the form:

$$X'(t) = g(X(t)), \quad X(0) = x_0, \quad t \in [0, 1], \quad (2)$$

where g is an unknown smooth function which is assumed to be positive on the range of $\{X(t) : t \in [0, 1]\}$. Therefore, the sample trajectory $X(t)$ is a strictly increasing function of time t . The observations are

$$Y_j = X(t_j) + \varepsilon_j, \quad j = 1, \dots, n, \quad (3)$$

where $0 \leq t_1, \dots, t_n \leq 1$ are observation times. The noise terms ε_j 's are assumed to be i.i.d. with mean 0 and variance $\sigma_\varepsilon^2 > 0$.

Our goal is to estimate the gradient function g based on the observed data Y_j s. We propose to approximate g through a basis representation:

$$g(x) \approx g_\beta := \sum_{k=1}^M \beta_k \phi_{k,M}(x), \quad (4)$$

where $\{\phi_{k,M}(\cdot)\}_{k=1}^M$ is a set of linearly independent compactly supported smooth functions. Henceforth, we use ϕ_k to denote $\phi_{k,M}$.

We now describe the estimation procedure. For the time being, assume that we observe the two endpoints $x_0 = X(0)$ and $x_1 = X(1)$ noiselessly and so the combined support of $\{\phi_1, \dots, \phi_M\}$ is the interval $[x_0, x_1]$. Given any $\beta := (\beta_1, \dots, \beta_M)$ so that g_β is positive on the support of $\{\phi_k(\cdot)\}_{k=1}^M$, we can solve the initial value problem

$$x'(t) = g_\beta(x(t)), \quad t \in [0, 1], \quad x(0) = x_0 \quad (5)$$

to obtain the corresponding trajectory $X(t; \beta)$. Define the L^2 loss function:

$$L(\beta) := \sum_{j=1}^n (Y_j - X(t_j; \beta))^2. \quad (6)$$

Then the proposed estimator of g is defined as

$$\hat{g}(x) := g_{\hat{\beta}}(x) = \sum_{k=1}^M \hat{\beta}_k \phi_{k,M}(x), \quad \text{where } \hat{\beta} := \arg \min_{\beta \in \mathbb{R}^M} L(\beta). \quad (7)$$

Minimization of $L(\beta)$ is a nonlinear least squares problem. We propose to use a Levenberg-Marquardt iterative updating scheme. Since this requires evaluating the trajectory $X(t; \beta)$ and its derivative with respect to β , given the current estimate of β , we solve the corresponding differential equations numerically by using the 4-th order Runge-Kutta method. More details are given in the Appendix. Finally, the number of basis M is selected through an approximate cross-validation score. A fitting procedure using similar techniques is studied in Paul et al. (2011) in a different context.

In practice, the initial value $x_0 = X(0)$ and the right boundary $x_1 = X(1)$ may not be observed or may be observed with noise. The choice of the endpoints of the combined support of the basis functions then becomes a delicate matter. This is because evaluation of the trajectory is an initial value problem, so error in x_0 propagates throughout the time domain. We discuss this in more details in the Appendix (particularly, see Figure A.1).

In the following, we propose a modified estimation procedure when x_0 and x_1 are unknown. The basic idea is to first estimate the trajectory at the endpoints of a slightly smaller time interval $[\delta, 1 - \delta]$ for a small positive constant δ , and then estimate the gradient function using data falling within this time interval. Throughout the paper, δ is treated as a fixed quantity. In practice, we may select δ to be the time point such that about 5% of the data fall in the intervals $[0, \delta]$ and $[1 - \delta, 1]$. Too small a value of δ may cause distortions of the estimated g at the boundaries.

We first obtain nonparametric estimates of $x_{0,\delta} := X(\delta)$ and $x_{1,\delta} := X(1 - \delta)$, denoted by \hat{x}_0 and \hat{x}_1 , respectively. We then define $x_{0,M} = \hat{x}_0 - \eta_M$ and $x_{1,M} = \hat{x}_1 + \eta_M$, where η_M is a small positive number satisfying $\eta_M = o(M^{-1})$ which implies that $x_0 < x_{0,M} < x_{1,M} < x_1$ as n goes to infinity. At the same time, η_M should be large enough so that $\max_{j=0,1} |x_{j,\delta} - \hat{x}_j| = o_P(\eta_M)$ which ensures that $x_{0,M} < x_{0,\delta} < x_{1,\delta} < x_{1,M}$ and $\max_{j=0,1} |x_{j,\delta} - x_{j,M}| = O_P(\eta_M) = o_P(M^{-1})$ as n goes to infinity. For some technical considerations, to be utilized later, we also want $\eta_M \gg M^{-3/2}$. In practice, we may select η_M to be $\min\{M^{-3/2} \log n, s_M / \log n\}$ where s_M is the length of the smallest support among the basis functions $\{\phi_1, \dots, \phi_M\}$. For more details on how to obtain \hat{x}_j , $j = 0, 1$, see Lemma 3.1 in Section 3. In addition, we also assume that \hat{x}_j , $j = 0, 1$ are estimated from a sample independent from that used in estimating β . This can be easily achieved in practice by sub-sampling of the measurements. This assumption enables us to prove the consistency result (in Section 3) conditionally on \hat{x}_j , $j = 0, 1$ and treating them as nonrandom sequences converging to $x_{j,\delta}$, $j = 0, 1$.

We then set the combined support of the basis functions $\{\phi_{k,M}\}_{k=1}^M$ as the interval $[x_{0,M}, x_{1,M}]$, and use the following modified loss function to derive an estimator for g :

$$\tilde{L}_\delta(\beta) = \sum_{j=1}^n (Y_j - X(t_j; \beta, \hat{x}_0))^2 \mathbf{1}_{[\delta, 1-\delta]}(t_j), \quad (8)$$

where $X(t; \beta, a)$ denotes the integral curve of the ODE

$$x'(t) = g_\beta(x(t)), \quad t \in [\delta, 1 - \delta], \quad x(\delta) = a. \quad (9)$$

The estimated \hat{g} is through minimizing the above loss function with respect to β (equation (7) with L replaced by \tilde{L}_δ).

3 Consistency

In this section, we discuss the consistency of the estimator \hat{g} defined by the loss function (8). The asymptotic framework is that the number of basis functions M goes to infinity together with the number of measure-

ments n . The consistency of the estimator \hat{g} over $[x_{0,\delta}, x_{1,\delta}]$ is formulated in terms of the L^2 -loss as:

$$\int_{x_{0,\delta}}^{x_{1,\delta}} |\hat{g}(u) - g(u)|^2 du \longrightarrow 0, \quad \text{in probability as } n \rightarrow \infty.$$

In Theorem 3.2 we derive a bound on the rate of convergence of \hat{g} in terms of the L^2 -loss as $n, M \rightarrow \infty$ that depends upon the degree of smoothness of g . Specifically, the optimal rate is $O_P(n^{-2p/(2p+3)})$ for $p \geq 4$.

3.1 Assumptions

The following assumptions are made on the model.

A1 $g \in C^p(D)$, and $g > 0$ on D for some integer $p \geq 3$, where D is an open interval containing $[x_0, x_1]$.

A2 The collection of basis functions $\Phi_M := \{\phi_{1,M}, \dots, \phi_{M,M}\}$ satisfies:

- (i) $\phi_{k,M}$'s have unit L^2 norm;
- (ii) the combined support of Φ_M is $D_0 \equiv D_{0,M} := [x_{0,M}, x_{1,M}]$ and for every k , the length of the support of $\phi_{k,M}$ is $O(M^{-1})$;
- (iii) $\phi_{k,M} \in C^2(D_0)$ for all k ;
- (iv) $\sup_{x \in D_0} \sum_{k=1}^M |\phi_{k,M}^{(j)}(x)|^2 = O(M^{1+2j})$, for $j = 0, 1, 2$;
- (v) the Gram matrix $\mathbf{G}_{\Phi_M} := ((\int_{x_{0,M}}^{x_{1,M}} \phi_{k,M}(u)\phi_{l,M}(u)du))_{k,l=1}^M$ is such that there exist constants $0 < \underline{c} \leq \bar{c} < \infty$, not depending on M such that $\underline{c} \leq \lambda_{\min}(\mathbf{G}_{\Phi_M}) \leq \lambda_{\max}(\mathbf{G}_{\Phi_M}) \leq \bar{c}$ for all M ;
- (vi) for every M , there is a $\beta^* \in \mathbb{R}^M$ such that $\sup_{t \in [\delta, 1-\delta]} |X_g(t) - X(t; \beta^*)| = O(M^{-(p+1)})$ and $\sup_{u \in [x_{0,\delta}, x_{1,\delta}]} |g^{(j)}(u) - g_{\beta^*}^{(j)}(u)| = O(M^{-p+j})$ for $j = 0, 1, 2$, where $g_{\beta} = \sum_{k=1}^M \beta_k \phi_{k,M}$ and $X(t; \beta) \equiv X(t; \beta, x_{0,\delta})$ with $X(t; \beta, a)$ as in (9).

A3 Time points $\{t_j\}_{j=1}^n$ are realizations of $\{T_j\}_{j=1}^n$, where T_j 's are i.i.d. from a continuous distribution F_T supported on $[0, 1]$ with a density f_T satisfying $\underline{c}' \leq f_T \leq \bar{c}'$ for some $0 < \underline{c}' \leq \bar{c}' < \infty$.

A4 The noise ε_j 's are i.i.d. sub-Gaussian random variables (cf. Vershynin, 2010) with mean 0 and variance $\sigma_\varepsilon^2 > 0$.

We give brief explanations of these assumptions. **A1** ensures sufficient smoothness of the solution paths of the differential equation (2). Also by **A1**, $X_g(\cdot)$ is $p+1$ times continuously differentiable on D . Assumptions (i) to (v) of **A2** are satisfied by B-spline basis, rescaled to have unit norm, of order ≥ 3 , with equally spaced knots. Define

$$\xi_n := \sqrt{\log nn}^{-\frac{p+1}{2p+3}}, \quad (10)$$

which is used in determining the rates of convergence of the estimator. Thus, by making use of **A4**, we get the following results with respect to the estimates of $x_{0,\delta}$ and $x_{1,\delta}$ (cf. Fan and Gijbels, 1996).

Lemma 3.1. *Suppose that **A1** and **A4** hold. Consider using a kernel of sufficient degree of smoothness to obtain estimates \hat{x}_j for $x_{j,\delta}$, $j = 1, 2$, through local polynomial method with bandwidth of order $n^{-1/(2p+3)}$. Define $d_n := \max_{j=0,1} |\hat{x}_j - x_{j,\delta}|$. Then $d_n = O_P(n^{-(p+1)/(2p+3)})$ and given $\eta > 0$, there exists $C(\eta) > 0$ such that $d_n \leq C(\eta)\xi_n$ with probability at least $1 - n^{-\eta}$, where ξ_n is as in (10).*

If $M = O((n/\log n)^{1/7})$ (as in Theorem 3.1) and $M^{-3/2} \ll \eta_M \ll M^{-1}$ for some $C > 0$, then we have $\xi_n = o(\eta_M)$ as $n \rightarrow \infty$. This ensures that $D_0 = [x_{0,M}, x_{1,M}]$ is within the interval $[x_0, x_1]$ a.s. for large enough n and hence the properties of the function g hold on D_0 . In addition, D_0 contains the interval $[x_{0,\delta}, x_{1,\delta}]$. Therefore condition (ii) in **A2** ensures that the combined support of the basis functions covers the range of the data used in estimating g .

Conditions (i) to (v) of **A2** are satisfied by many classes of basis functions, including normalized B-spline basis of order ≥ 3 with equally spaced knots in the interval $[x_{0,M}, x_{1,M}]$. We show in Appendix B that if the B-splines basis of order $\geq \max\{3, p-1\}$ with equally spaced knots in $[x_{0,M}, x_{1,M}]$, then (vi) of **A2** is also satisfied. Condition (vi) of **A2** ensures that a solution $X(t; \beta)$ of (5) on $t \in [\delta, 1-\delta]$ exists for all β sufficiently close to β^* . This allows us to apply the perturbation theory of differential equations to bound the fluctuations of the sample paths when we perturb the parameter β .

Assumption **A3** on the randomness of the sample points allows us to work with the random variables \tilde{T}_j defined as T_j conditional on $T_j \in [\delta, 1-\delta]$ with conditional density \tilde{f}_T given by $\tilde{f}_T(t) = f_T(t)/(F_T(1-\delta) - F_T(\delta))$. The properties of f_T ensure that \tilde{f}_T satisfies the same property on $[\delta, 1-\delta]$ with possibly modified values of the constants c_1 and c_2 . It should be noted that the key derivations leading to the consistency of \hat{g} are conditional on \mathbf{T} and therefore **A3** is only a convenient assumption for describing the regularity of the time points. The asymptotic results (Theorems 3.1 and 3.2) hold if instead of being randomly distributed, the time points form a fixed regular grid, say, with equal spacing.

3.2 Rate of convergence

As mentioned earlier, the estimation of $g(\cdot)$ is a nonlinear inverse problem since $X'(t)$ is not directly observable. In addition, this is also an ill-posed estimation problem. Let $X^\beta(\cdot; \beta)$ be the partial derivative of $X(\cdot; \beta)$ with respect to β , where $X(\cdot; \beta) \equiv X(\cdot; \beta, x_{0,\delta})$ is the solution of (9) with $x(0) = x_{0,\delta}$. Let $\beta^* \in \mathbb{R}^M$ be as in **A2**. Define

$$G_* := \mathbb{E} \left(X^\beta(\tilde{T}_1; \beta^*) (X^\beta(\tilde{T}_1; \beta^*))^T \right), \quad (11)$$

where the expectation is with respect to the distribution of \tilde{T}_1 . Clearly G_* is a positive semi-definite matrix. It becomes clear from the analysis carried out later that the degree of ill-posedness of the estimation problem is determined by the size of the operator norm of the matrix G_* as a function of M . The following proposition gives a precise quantification of the degree of ill-posedness. The situation here is in contrast with standard nonparametric function estimation problems where the corresponding matrix is well-conditioned.

Proposition 3.1. *Assume that assumptions **A1** to **A3** hold with $p \geq 3$. Assume further that (a) $\max_{j=0,1} |x_{j,M} - x_{j,\delta}| = o(M^{-1})$ (a.s.) and (b) $\min\{x_{1,M} - x_{1,\delta}, x_{0,\delta} - x_{0,M}\} \gg M^{-3/2}$. Then (a.s.)*

$$\|G_*^{-1}\| = O(M^2). \quad (12)$$

By Lemma 3.1 and the discussion that follows, under the condition of Theorem 3.1, (a) and (b) of Proposition 3.1 hold.

We now state the main result on the consistency of the estimate \hat{g} .

Theorem 3.1. *Suppose that the observed data $\{Y_j : j = 1, \dots, n\}$ follow the model described by equations (2) and (3) and that assumptions **A1**–**A4** are satisfied with $p \geq 3$. Suppose further that the sequence M is such that*

$$c'_1 \left(\frac{n}{\sigma_\varepsilon^2} \right)^{1/(2p+3)} \leq M \ll \left(\frac{n}{\sigma_\varepsilon^2 \log n} \right)^{1/7} \quad (13)$$

for some $c'_1 > 0$, $M^{-3/2} \ll \eta_M \ll M^{-1}$, and ξ_n be as defined in Lemma 3.1. Let $\bar{\alpha}_n := c'_2 M^{-2}$ for some $c'_2 > 0$ (sufficiently small) and

$$\alpha_n := C_0 M \max \left\{ \sigma_\varepsilon \sqrt{\frac{M \log n}{n}}, M^{-(p+1)}, \xi_n \right\}, \quad (14)$$

for some $C_0 > 0$. Then as $n \rightarrow \infty$, with probability tending to one, there exists a local minimum $\hat{\beta}$ of the objective function $\tilde{L}_\delta(\beta)$ (defined through (8)), which is also a global minimum within radius $\bar{\alpha}_n$ of β^* (note that, $\alpha_n \leq \bar{\alpha}_n$ by (13)) such that, with $\hat{g} := g_{\hat{\beta}}$,

$$\int_{x_{0,\delta}}^{x_{1,\delta}} |\hat{g}(u) - g(u)|^2 du = O(\alpha_n^2). \quad (15)$$

The proof of Theorem 3.1 is given in Section 4.

Remark 3.1. Assuming σ_ε to be a constant, if M is chosen to be of the order $n^{1/(2p+3)}$, then α_n^2 in (15) simplifies to $n^{-2p/(2p+3)} \log n$, which is within a factor of $\log n$ of the optimal rate in terms of the L^2 -loss for estimating $X'(t)$ based on the data $\{Y_j : j = 1, \dots, n\}$ given by (2) when $X \in C^{p+1}([0, 1])$. The fact that an estimator of g can attain this rate can be anticipated from the representation of g as $g = X' \circ X^{-1}$. For $p \geq 4$, we can improve the rate of convergence of \hat{g} slightly further, by dropping the factor of $\log n$, as stated in the following result.

Theorem 3.2. Suppose that the conditions of Theorem 3.1 are satisfied with $p \geq 4$ and, further, the sequence M satisfies the condition that $c(n/\sigma_\varepsilon^2)^{1/(2p+3)} \leq M \ll (n/\sigma_\varepsilon^2 \log n)^{1/9}$ for some $c > 0$. Let \hat{g} be as in Theorem 3.1. Then,

$$\begin{aligned} & \int_{x_{0,\delta}}^{x_{1,\delta}} (\hat{g}(x) - g(x))^2 dx \\ &= O_P \left(\frac{\sigma_\varepsilon^2 M^3}{n} \right) + O_P(M^{-2p}) + O_P \left(M^2 (\sigma_\varepsilon^2/n)^{2(p+1)/(2p+3)} \right), \end{aligned} \quad (16)$$

with the optimal rate given by $O_P((\sigma_\varepsilon^2/n)^{2p/(2p+3)})$, which is obtained when $M = c(n/\sigma_\varepsilon^2)^{1/(2p+3)}$ for some $c > 0$.

Proof of Theorem 3.2 is given in Section S2 of SM.

We can also derive an approximate expression for the asymptotic variance of $\hat{g}(\cdot)$. Using a consistent root $\hat{\beta}$, we can use the equation $\nabla L(\beta)|_{\beta=\hat{\beta}} = 0$. Using the asymptotic representation of $\hat{\beta} - \beta^*$ used in the proof of Theorem 3.2 (see Section S2 of SM), and ignoring higher order terms and the contribution of the model bias, and finally evaluating the expressions at β instead of β^* (which is unknown), we have

$$\text{Var}(\hat{\beta}) \approx D(\hat{\beta}) := \hat{\sigma}_\varepsilon^2 \left[\sum_{j=1}^n \left(\frac{\partial X(T_j; \hat{\beta})}{\partial \beta} \right) \left(\frac{\partial X(T_j; \hat{\beta})}{\partial \beta} \right)^T \right]^{-1}. \quad (17)$$

Here the estimated noise variance $\hat{\sigma}_\varepsilon^2$ can be computed as the mean squared error $(n - M)^{-1} \sum_{j=1}^M (Y_j - X(T_j; \hat{\beta}))^2$. The expression (17) allows us to obtain an approximate asymptotic variance for $\hat{g}(x)$ by $V(x) := \phi(x)^T D(\hat{\beta}) \phi(x)$, for any given x , where $\phi(x) = (\phi_{1,M}(x), \dots, \phi_{M,M}(x))^T$.

3.3 Initial estimator

In Theorem 3.1 we prove the rate of convergence for a local minimizer, which is a global minimizer within a radius of $O(M^{-2})$ of β^* for a suitable range of values of M . Therefore, we need an initial estimate which resides within this domain. In the following, we describe one way of obtaining such an initial estimate, through a two-stage approach, which is similar in spirit to the approaches by Chen and Wu (2008a, 2008b).

Suppose that we first estimate $X(t)$ and $X'(t)$ by local polynomial smoothing and denote these estimates by $\hat{X}(t)$ and $\hat{X}'(t)$. Then, we fit the regression model

$$\hat{X}'(T_j) = \phi(\hat{X}(T_j))^T \beta + e_j, \quad j = 1, \dots, n \quad (18)$$

by ordinary least squares, where $\phi = (\phi_{1,M}, \dots, \phi_{M,M})$. We refer to the resulting estimator $\tilde{\beta}$ as the two-stage estimator of β :

$$\tilde{\beta} = \left[\sum_{j=1}^n \phi(\hat{X}(T_j)) \phi(\hat{X}(T_j))^T \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right]^{-1} \left(\sum_{j=1}^n \hat{X}'(T_j) \phi(\hat{X}(T_j)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right). \quad (19)$$

Since $X(t)$ is $p+1$ times continuously differentiable and $X'(t)$ is p -times continuously differentiable (by **A1**), and $\{\varepsilon_j\}$ is sub-Gaussian, with the optimal choice of bandwidths, we have

$$\max_{1 \leq j \leq n} |\hat{X}(T_j) - X(T_j)| \mathbf{1}_{[\delta, 1-\delta]}(T_j) = O((\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)} \sqrt{\log n}) \quad (20)$$

$$\max_{1 \leq j \leq n} |\hat{X}'(T_j) - X'(T_j)| \mathbf{1}_{[\delta, 1-\delta]}(T_j) = O((\sigma_\varepsilon^2/n)^{p/(2p+3)} \sqrt{\log n}) \quad (21)$$

with probability tending to 1. We state the following result about the rate of convergence of the two-stage estimator. The proof is given in Section S3 of SM.

Proposition 3.2. *Suppose that $p \geq 2$ and **A1–A4** hold and that the two-stage estimate of g is given by $\tilde{g} := g_{\tilde{\beta}}$ where $\tilde{\beta}$ is defined in (19). Then, supposing that $n^{1/(4p+6)} \ll M \ll n^{(p+1)/(4p+6)} / \sqrt{\log n}$, with probability tending to 1,*

$$\int_{x_{0,\delta}}^{x_{1,\delta}} |\tilde{g}(u) - g(u)|^2 du = O(\tilde{\alpha}_n^2) \quad (22)$$

where

$$\tilde{\alpha}_n = \max\{M^2(\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)} \sqrt{\log n}, M^{-p}\}. \quad (23)$$

When $\sigma_\varepsilon \asymp 1$, the optimal value of $\tilde{\alpha}_n$ is of the order $n^{-p(p+1)/(p+2)(2p+3)} (\log n)^{-p/(2p+4)}$ is obtained when $M \asymp M_* = n^{(p+1)/(p+2)(2p+3)} (\log n)^{-1/(2p+4)}$. It can be checked that for all $p \geq 3$, this rate is slower than the optimal α_n for the nonlinear regression-based estimator \hat{g} derived in Theorem 3.1. However, the rate of convergence of this estimator is faster than $O(M^{-2}) = O(\bar{\alpha}_n)$ if $M^4 \ll n^{(p+1)/(2p+3)} / \sqrt{\log n}$. So, for these range of M , which includes M_* , the two-stage estimator resides within the ball of radius $O(\bar{\alpha}_n)$ around β^* , over which \hat{g} , the optimizer of (8), is a global optimum.

4 Proofs

In this section, we outline the main steps of the proof. Some technical details are deferred to the Appendix.

The main idea behind the proof of Theorem 3.1 is to obtain a lower bound on the difference $n^{-1}(L_\delta(\beta) - L_\delta(\beta^*))$ which is proportional $\|\beta - \beta^*\|^2$ when β lies in an annular region around β^* . The outer radius of the annular region depends on the degree of ill-conditioning of the problem, as quantified by Proposition 3.1, and the smoothness of the function g and the approximating bases, as indicated in condition **A2**. This lower bound then naturally leads to the conclusion about the existence and rate of convergence of a local minimizer \hat{g} .

Proof of Proposition 3.1

For convenience of notations, we define $X_*(t)$ to be the sample path $X(t; \beta^*)$. Since $X^\beta(\cdot; \beta)$ is given by (A.5) in the Appendix,

$$X^{\beta_r}(t) = g_{\beta}(X(t)) \int_{x_0}^{X(t)} \frac{\phi_r(x)}{(g_{\beta}(x))^2} dx, \quad r = 1, \dots, M,$$

in order to prove Proposition 3.1, it suffices to find a lower bound on

$$\min_{\|\mathbf{b}\|=1} \int_{\delta}^{1-\delta} \left[\int_{\delta}^t g_{\mathbf{b}}(X_*(u)) / g_{\beta^*}(X_*(u)) du \right]^2 \tilde{f}_T(t) dt$$

where $g_{\mathbf{b}}(u) = \mathbf{b}^T \phi(u)$ with $\phi = (\phi_1, \dots, \phi_M)^T$. By **A3**, without loss of generality, we can take the density $\tilde{f}_T(\cdot)$ to be uniform on $[\delta, 1 - \delta]$.

We make use of the following result known as Halperin-Pitt inequality (Mitrinovic *et al.*, 1991).

Lemma 4.1. *If f is locally absolutely continuous and f'' is in $L_2([0, A])$, then for any $\epsilon > 0$ the following inequality holds with $K(\epsilon) = 1/\epsilon + 12/A^2$,*

$$\int_0^A (f'(t))^2 dt \leq K(\epsilon) \int_0^A f^2(t) dt + \epsilon \int_0^A (f''(t))^2 dt. \quad (24)$$

Now defining

$$R(t) := \int_{\delta}^t \frac{g_{\mathbf{b}}(X_*(u))}{g_{\beta^*}(X_*(u))} du,$$

we have,

$$\begin{aligned} R'(t) &:= \frac{dR(t)}{dt} = \frac{g_{\mathbf{b}}(X_*(t))}{g_{\beta^*}(X_*(t))} \\ R''(t) &:= \frac{d^2 R(t)}{dt^2} = \left[\frac{g'_{\mathbf{b}}(X_*(t))}{g_{\beta^*}(X_*(t))} - \frac{g_{\mathbf{b}}(X_*(t)) g'_{\beta^*}(X_*(t))}{g_{\beta^*}^2(X_*(t))} \right] X'_*(t) \\ &= \left[\frac{g'_{\mathbf{b}}(X_*(t))}{g_{\beta^*}(X_*(t))} - \frac{g_{\mathbf{b}}(X_*(t)) g'_{\beta^*}(X_*(t))}{g_{\beta^*}^2(X_*(t))} \right] g_{\beta^*}(X_*(t)). \end{aligned}$$

By (vi) of **A2**, we have $\sup_{t \in [\delta, 1-\delta]} |X_g(t) - X_*(t)| = O(M^{-(p+1)})$ and hence

$$X_*(1 - \delta) \leq x_{1,\delta} + |X_*(1 - \delta) - x_{1,\delta}| < x_{1,M}, \quad X_*(\delta) \geq x_{0,\delta} - |X_*(\delta) - x_{0,\delta}| > x_{0,M}. \quad (25)$$

Hence, using the facts that the coordinates of $\phi(u)$ are $O(M^{1/2})$ and the coordinates of $\phi'(u)$ are $O(M^{3/2})$, and all these functions are supported on intervals of length $O(M^{-1})$, we deduce that,

$$\int_{\delta}^{1-\delta} (R''(t))^2 dt = O(M^2). \quad (26)$$

An application of Lemma 4.1 with $f(t) = R(t - \delta)$ and $A = 1 - 2\delta$ yields

$$\int_{\delta}^{1-\delta} (R'(t))^2 dt \leq (1/\epsilon + 12/(1 - 2\delta)^2) \int_{\delta}^{1-\delta} (R(t))^2 dt + \epsilon \int_{\delta}^{1-\delta} (R''(t))^2 dt. \quad (27)$$

Take $\epsilon = k_0 M^{-2}$ for some $k_0 > 0$, then by (26),

$$\int_{\delta}^{1-\delta} (R(t))^2 dt \geq k_1 M^{-2} \int_{\delta}^{1-\delta} (R'(t))^2 dt - k_2 M^{-2},$$

for constants $k_1, k_2 > 0$ dependent on k_0 . Next, we write

$$\int_{\delta}^{1-\delta} (R'(t))^2 dt = \int_{X_*(\delta)}^{X_*(1-\delta)} \frac{g_{\mathbf{b}}^2(v)}{g_{\beta^*}^3(v)} dv = \int_{X_*(\delta)}^{X_*(1-\delta)} g_{\mathbf{b}}^2(v) h(v) dv \quad (28)$$

where $h(v) = g_{\beta^*}^{-3}(v)$ which is bounded below by a positive constant on the interval $[X_*(\delta), X_*(1-\delta)]$.

Observe that by (25), the combined support of $\{\phi_{k,M}\}_{k=1}^M$, viz., $[x_{0,M}, x_{1,M}]$, contains (for sufficiently large M) the interval $[X_*(\delta), X_*(1-\delta)]$. Also, $|x_{1,M} - X_*(1-\delta)| \leq |x_{1,M} - x_{1,\delta}| + |x_{1,\delta} - X_*(1-\delta)| = o(M^{-1})$ and $|x_{0,M} - X_*(\delta)| \leq |x_{0,M} - x_{0,\delta}| + |x_{0,\delta} - X_*(\delta)| = o(M^{-1})$. These two facts and the condition (v) of **A2** imply that

$$\begin{aligned} & \int_{X_*(\delta)}^{X_*(1-\delta)} g_{\mathbf{b}}^2(v) h(v) dv \\ & \geq \left(\inf_{v \in [X_*(\delta), X_*(1-\delta)]} h(v) \right) \mathbf{b}^T \left[\int_{x_{0,M}}^{x_{1,M}} \phi(v) (\phi(v))^T dv - o(1) \right] \mathbf{b} \geq k_3, \end{aligned}$$

for some constant $k_3 > 0$, for sufficiently large M . Thus, by appropriate choice of ϵ , we have $\int_{\delta}^{1-\delta} (R(t))^2 dt \geq k_4 M^{-2}$ for some constant $k_4 > 0$, which yields (12).

Proof of Theorem 3.1

Define

$$\Gamma_n(\beta, \beta^*) = \frac{1}{n} \sum_{j=1}^n (X(T_j; \beta) - X(T_j; \beta^*))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j), \quad (29)$$

$\mathcal{A}_M(\alpha_n, \bar{\alpha}_n) = \{\beta \in \mathbb{R}^M : \alpha_n \leq \|\beta - \beta^*\| \leq \bar{\alpha}_n\}$, and

$$D_n^* = \frac{1}{n} \sum_{j=1}^n (X_g(T_j) - X(T_j; \beta^*))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j).$$

Suppose that $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$. Henceforth, we use $X(t; \beta)$ to denote $X(t; \beta; x_{0,\delta})$ and $X_g(t)$ to denote $X_g(t; x_{0,\delta})$. Then

$$\begin{aligned} & \frac{1}{n} L_{\delta}(\beta) - \frac{1}{n} L_{\delta}(\beta^*) \\ &= \frac{1}{n} \sum_{j=1}^n (Y_j - X(T_j; \beta))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j) - \frac{1}{n} \sum_{j=1}^n (Y_j - X(T_j; \beta^*))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ &= \frac{1}{n} \sum_{j=1}^n (X(T_j; \beta) - X(T_j; \beta^*))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ & \quad - \frac{2}{n} \sum_{j=1}^n \varepsilon_j (X(T_j; \beta) - X(T_j; \beta^*)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ & \quad - \frac{2}{n} \sum_{j=1}^n (X_g(T_j) - X(T_j; \beta^*)) (X(T_j; \beta) - X(T_j; \beta^*)) \mathbf{1}_{[\delta, 1-\delta]}(T_j), \end{aligned} \quad (30)$$

where $U_{1n}(\beta, \beta^*)$ and $U_{2n}(\beta, \beta^*)$, are the second and third summations in the above expression, respectively. Next, we write

$$\begin{aligned}
& \frac{1}{n} \tilde{L}_\delta(\beta) - \frac{1}{n} L_\delta(\beta) \\
&= \frac{1}{n} \sum_{j=1}^n [(Y_j - X(T_j; \beta, \hat{x}_0))^2 - (Y_j - X(T_j; \beta))^2] \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&= \frac{1}{n} \sum_{j=1}^n (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad - \frac{2}{n} \sum_{j=1}^n \varepsilon_j (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad - \frac{2}{n} \sum_{j=1}^n (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta))(X_g(T_j) - X(T_j; \beta)) \mathbf{1}_{[\delta, 1-\delta]}(T_j), \tag{31}
\end{aligned}$$

where $V_{1n}(\beta)$, $V_{2n}(\beta)$, $V_{3n}(\beta)$ are the three summations in the last expression.

From (30) and (31), we deduce that

$$\begin{aligned}
& \frac{1}{n} \tilde{L}_\delta(\beta) - \frac{1}{n} \tilde{L}_\delta(\beta^*) \\
&= \frac{1}{n} (L_\delta(\beta) - L_\delta(\beta^*)) + \frac{1}{n} (\tilde{L}_\delta(\beta) - L_\delta(\beta)) - \frac{1}{n} (\tilde{L}_\delta(\beta^*) - L_\delta(\beta^*)) \\
&= \Gamma_n(\beta, \beta^*) - U_{1n}(\beta, \beta^*) - U_{2n}(\beta, \beta^*) + (V_{1n}(\beta) - V_{1n}(\beta^*)) \\
&\quad - (V_{2n}(\beta) - V_{2n}(\beta^*)) + U_{3n}(\beta, \beta^*) - U_{4n}(\beta, \beta^*) \tag{32}
\end{aligned}$$

where

$$\begin{aligned}
U_{3n}(\beta, \beta^*) &= \frac{2}{n} \sum_{j=1}^n (X_g(T_j) - X(T_j; \beta^*)) (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad - \frac{2}{n} \sum_{j=1}^n (X_g(T_j) - X(T_j; \beta^*)) (X(T_j; \beta^*; \hat{x}_0) - X(T_j; \beta^*)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
U_{4n}(\beta, \beta^*) &= \frac{2}{n} \sum_{j=1}^n (X(T_j; \beta) - X(T_j; \beta^*)) (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta)) \mathbf{1}_{[\delta, 1-\delta]}(T_j).
\end{aligned}$$

Using the fact that $\equiv X^a(t; \beta, a_0) := (\partial/\partial a)X(t; \beta, a)|_{a=a_0}$ can be expressed as

$$X^a(t; \beta, a_0) = \frac{g_\beta(X(t; \beta, a_0))}{g_\beta(a_0)}, \quad t \in [\delta, 1-\delta],$$

provided $g_\beta(x) > 0$ for $x \in [a_0, X(1-\delta; \beta, a_0)]$, we have, for all $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$,

$$\sup_{a_0 \in [x_{0,\delta} - \xi_n, x_{0,\delta} + \xi_n]} \sup_{t \in [\delta, 1-\delta]} |X(t; \beta, a_0) - X(t; \beta, x_{0,\delta})| \leq C_1 \xi_n \tag{33}$$

for some $C_1 > 0$. Here, we have used the fact that for $t \in [\delta, 1-\delta]$, and $a_0 \in [x_{0,\delta} - \xi_n, x_{0,\delta} + \xi_n]$,

$$X(t; \beta, a_0) = \tilde{G}_\beta^{-1}(t - \delta + G_\beta(a_0)) \quad \text{where} \quad \tilde{G}_\beta(x) := \int_{x_{0,M}}^x \frac{du}{g_\beta(u)}, \tag{34}$$

and that

$$\sup_{\beta \in \mathcal{A}(\alpha_n, \bar{\alpha}_n)} \sup_{x \in [x_{0,M}, x_{1,M}]} |g_\beta(x) - g_{\beta^*}(x)| = O(\bar{\alpha}_n M^{1/2}) = O(M^{-3/2}),$$

so that, by using (A.11), and the fact that $M^{-3/2} \ll \eta_M \ll M^{-1}$,

$$[x_{0,\delta} - \xi_n, \sup_{\beta \in \mathcal{A}(\alpha_n, \bar{\alpha}_n)} \sup_{a_0 \in [x_{0,\delta} - \xi_n, x_{0,\delta} + \xi_n]} X(1 - \delta; \beta, a_0)] \subset [x_{0,M}, x_{1,M}]$$

for large enough M and n .

We now bound individual terms in the expansion (32). First, we have the following lower bound on $\Gamma_n(\beta, \beta^*)$, the proof of which is given in Appendix C.

Lemma 4.2. *Let $\Gamma_n(\beta, \beta^*)$ be as defined in (29). Then given $\eta > 0$, there exist constants $d_1(\eta) > 0$ and $d_2, d_3 > 0$ independent of η such that*

$$\Gamma_n(\beta, \beta^*) \geq d_1(\eta) \frac{1}{M^2} \|\beta - \beta^*\|^2 - d_2 \|\beta - \beta^*\|^4 M^2 (1 + d_3 \bar{\alpha}_n^2 M^3) \quad (35)$$

uniformly in $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$ with probability at least $1 - n^{-\eta}$.

Since $\bar{\alpha}_n M^{3/2} = c'_2 M^{-1/2} = o(1)$, and the constant c'_2 can be chosen to be small enough so that we can conclude from (35) that given $\eta > 0$, there exists $d_4(\eta) > 0$ such that

$$\mathbb{P} \left(\Gamma_n(\beta, \beta^*) \geq \frac{d_4(\eta)}{M^2} \|\beta - \beta^*\|^2 \text{ for all } \beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n) \right) \geq 1 - n^{-\eta}. \quad (36)$$

Next, by Cauchy-Schwarz inequality, we have

$$|U_{2n}(\beta, \beta^*)| \leq 2\sqrt{D_n^*} \sqrt{\Gamma_n(\beta, \beta^*)}. \quad (37)$$

Next, by (33), we have

$$\max\{V_{1n}(\beta^*), \sup_{\beta \in \mathcal{A}(\alpha_n, \bar{\alpha}_n)} V_{1n}(\beta)\} \leq C_1^2 \xi_n^2, \quad (38)$$

and hence

$$\sup_{\beta \in \mathcal{A}(\alpha_n, \bar{\alpha}_n)} |U_{3n}(\beta, \beta^*)| \leq 4C_1 \xi_n \sqrt{D_n^*} \quad (39)$$

and

$$|U_{4n}(\beta, \beta^*)| \leq 2C_1 \xi_n \sqrt{\Gamma_n(\beta, \beta^*)}. \quad (40)$$

Next, defining

$$Z(\beta) = \frac{\sum_{j=1}^n \varepsilon_j (X(T_j; \beta) - X(T_j; \beta^*)) \mathbf{1}_{[\delta, 1-\delta]}(T_j)}{\sigma_\varepsilon \sqrt{\sum_{j=1}^n (X(T_j; \beta) - X(T_j; \beta^*))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j)}},$$

and setting $Z(\beta)$ being zero if the denominator is zero, we have

$$|U_{1n}(\beta)| \leq \frac{2\sigma_\varepsilon}{\sqrt{n}} \sqrt{\Gamma_n(\beta, \beta^*)} |Z(\beta)|. \quad (41)$$

Let $\mathcal{B}_M(\Delta; \alpha_n, \bar{\alpha}_n)$ be a Δ -net for $\mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$. Then $|\mathcal{B}_M(\Delta; \alpha_n, \bar{\alpha}_n)| \leq 3(\bar{\alpha}_n/\Delta)^M$. Then, by using Lemma S.3 in SM, and (36), we conclude that given $\eta > 0$, there exist constants $c_1(\eta) > 0$, $C'(\eta) > 0$, and a set $A_{1\eta}$ with $\mathbb{P}(\mathbf{T} \in A_{1\eta}) \geq 1 - n^{-\eta}$, such that for all $\mathbf{T} \in A_{1\eta}$,

$$\mathbb{P} \left(\max_{\beta \in \mathcal{B}_M(\delta; \alpha_n, \bar{\alpha}_n)} |Z(\beta)| > c_1(\eta) \sqrt{M \log(\bar{\alpha}_n/\delta)} \mid \mathbf{T} \right) \leq C'(\eta) \left(\frac{\Delta}{\bar{\alpha}_n} \right)^{\eta M}$$

for some constant $C' > 0$. Thus, taking δ to be sufficiently small, say, $\delta = n^{-c}$ for c large enough, and using the smoothness of the process $Z(\beta)$ as a function of β , we can show that given any $\eta > 0$, there exists $c_2(\eta) > 0$, such that for all $\mathbf{T} \in A_{1\eta}$,

$$\mathbb{P} \left(\sup_{\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)} |Z(\beta)| \leq c_2(\eta) \sqrt{M \log n} \mid \mathbf{T} \right) > 1 - n^{-\eta}. \quad (42)$$

Very similarly, defining

$$\tilde{Z}(\beta) = \frac{\sum_{j=1}^n \varepsilon_j (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta)) \mathbf{1}_{[\delta, 1-\delta]}(T_j)}{\sigma_\varepsilon \sqrt{\sum_{j=1}^n (X(T_j; \beta; \hat{x}_0) - X(T_j; \beta))^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j)}},$$

expressing $V_{2n}(\beta) = 2\sigma_\varepsilon n^{-1/2} \sqrt{V_{1n}(\beta)} \tilde{Z}(\beta)$, and using (38), we have, for any given $\eta > 0$, there exists $c_3(\eta) > 0$ and a set $A_{2\eta}$ with $\mathbb{P}(\mathbf{T} \in A_{2\eta}) \geq 1 - n^{-\eta}$, such that for all $\mathbf{T} \in A_{2\eta}$,

$$\mathbb{P} \left(\sup_{\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n) \cup \{\beta^*\}} |V_{2n}(\beta)| \leq c_3(\eta) \sigma_\varepsilon \xi_n \sqrt{\frac{M \log n}{n}} \mid \mathbf{T} \right) > 1 - n^{-\eta}. \quad (43)$$

Finally, by **A2** we have the bound

$$D_n^* \leq \sup_{t \in [\delta, 1-\delta]} |X_g(t) - X(t; \beta^*)|^2 \leq C_2 M^{-2(p+1)} \quad (44)$$

for some $C_2 > 0$.

Combining (37)–(44), we claim that, given $\eta > 0$, there exist constants $C_3(\eta) > 0$, $C_4(\eta) > 0$, and constants $C_l > 0$, $l = 5, \dots, 8$, not depending on η , such that uniformly on $\mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$

$$\begin{aligned} & \frac{1}{n} \tilde{L}_\delta(\beta) - \frac{1}{n} \tilde{L}_\delta(\beta^*) \\ & \geq \Gamma_n(\beta, \beta^*) - \sqrt{\Gamma_n(\beta, \beta^*)} \left(C_3(\eta) \sqrt{\frac{M \log n}{n}} + C_5 M^{-(p+1)} + C_6 \xi_n \right) \\ & \quad - \xi_n \left(C_4(\eta) \sqrt{\frac{M \log n}{n}} + C_7 M^{-(p+1)} + C_8 \xi_n \right) \end{aligned} \quad (45)$$

with probability at least $1 - O(n^{-\eta})$.

From (45) and (36), and a careful choice of the constant C_0 in the definition (14) of α_n , and with M as in 13, we conclude that for any $\eta > 0$, there exists $C_9(\eta) > 0$ such that, uniformly in $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$,

$$\frac{1}{n} \tilde{L}_\delta(\beta) - \frac{1}{n} \tilde{L}_\delta(\beta^*) \geq C_9(\eta) \frac{1}{M^2} \|\beta - \beta^*\|^2 \quad (46)$$

with probability at least $1 - O(n^{-\eta})$. From this, we can conclude that with probability at least $1 - O(n^{-\eta})$ there exists a local minimum $\hat{\beta}$ of $\tilde{L}_\delta(\beta)$, which is also a global minimum within radius $\bar{\alpha}_n$ of β^* and which satisfies $\|\hat{\beta} - \beta^*\| = O(\alpha_n)$.

5 Simulation Study

In this section, we conduct a simulation study to examine the finite sample performance of the proposed estimation procedure, as well as to compare it with the two-stage estimator described in Section 3.3.

In the simulation, the true gradient function g is represented by 4 B-spline functions with knots at 0.35, 0.60, 0.85, 1.10 and respective coefficients 0.1, 1.2, 1.6, 0.4 (shown by the blue curve in Figure 1). We set the initial value $X(0) = x_0 = 0.25$ in equation (2) to generate the true trajectory $X(\cdot)$. We then simulate 100 independent data sets according to equation (3). Specifically, for each data set, we first randomly choose an integer n from $\{60, \dots, 100\}$. Then n observation times $\{t_1, \dots, t_n\}$ are uniformly sampled from $[0, 1]$. Finally, the Y_j 's are generated according to equation (3) with added noise $\epsilon_i \sim \text{Normal}(0, 0.01^2)$. The observed data from one such replicate is shown in Figure S.1 in SM together with the true trajectory $X(\cdot)$.

We fit the proposed estimator $\hat{g}(\cdot)$ with M B-spline basis functions with equally spaced knots on $[0.1, 1.1]$. We consider $M = 3, 4, 5$ and choose M by an approximate leave-one-out CV score criterion similar to that used in Paul et al. (2011). Out of the 100 replicates, 43 times the model with $M = 4$ (the true model) is chosen and 66 times the model with $M = 5$ is chosen.

We also consider the two-stage estimator, where in the first stage, the sample trajectory $X(\cdot)$ and its derivative $X'(\cdot)$ are estimated by applying local linear and local quadratic smoothing with Gaussian Kernel, respectively, to the observed data $\{(t_j, Y_j)\}_{j=1}^n$. The bandwidths are chosen by cross-validation. In the second stage, a quadratic smoothing of $\hat{X}'(\cdot)$ versus $\hat{X}(\cdot)$ is performed to get an estimate of $g(\cdot)$.

Figure 1 shows the estimated gradient functions (red curves) of these 100 independent replicates overlaid on the true gradient function (blue curve). It can be seen from this figure that, the proposed estimator shows little bias. Its sampling variability is somewhat larger on the left side of the observed x domain than on the right side of the observed x domain. It performs much better than the two-stage estimator which shows both high bias and high variance. Indeed, the bias of the two-stage estimator would not go away even when in the second stage the true model is used to estimate g (through a least-squares regression of $\hat{X}'(\cdot)$ versus $\hat{X}(\cdot)$).

Figure S.2 shows the estimated trajectories (red curves) of these 100 independent replicates overlaid on the true trajectory (blue curve). In the left panel of the figure, the estimated trajectories are solved from equation (2) using the 4th-order Runge-Kutta method with g being the proposed estimator $\hat{g}(\cdot)$. In the right panel of the figure, the trajectories are estimated by applying local linear smoothing of the observed data (which are then used in the two-stage fitting for $g(\cdot)$). The estimated trajectories from the proposed procedure follow the true trajectory very well with little bias, whereas the estimator from the first-stage smoothing of the two-stage procedure shows more bias and more variability. Figure S.3 in SM shows the estimated derivative of the trajectory. Again, the proposed procedure gives a much better estimate of $X'(\cdot)$ than the presmoothing estimate (by local quadratic smoothing) used in the two-stage procedure.

6 Application : Berkeley Growth Data

In this section, we apply the proposed model to the Berkeley growth data (Tuddenham and Snyder, 1954). Although in the literature, there are many studies of growth curves (Hauspie et al., 1980; Milani, 2000), most of them try to model either the growth trajectories (i.e., $X(\cdot)$) or the rate of growth (i.e., $X'(\cdot)$). On the contrary, our goal is to estimate the gradient function, i.e., the functional relationship between $X'(\cdot)$ and $X(\cdot)$ which provides insights of the growth dynamics, such as at what height the growth rate tends to be the highest.

Specifically, we fit the proposed model to each of the 54 female subjects in this data set. For each girl, her heights were measured at 31 time points from 1 year old to 18 years old. We use M B-spline basis

functions with equally spaced knots. We consider $M = 4, 5, 6, 7$ and for each subject we choose the “best” M using an approximate leave-one-out CV score. In 37 out of 54 subjects, the model with $M = 6$ is chosen, and for the rest 17 subjects, the model with $M = 7$ is chosen. Figure 2 shows the fitted gradient functions for these 54 subjects. From this figure, we can see that, most girls experienced two growth spurs, one at the birth (when their heights are shortest) and another when they were around either 130 cm tall or 150 cm tall. Moreover Figure S.4 in SM shows the fitted gradient functions with the two-standard-error bands (by equation (17)) for 25 girls. Figure S.5 in SM shows the observed (red dots) and fitted (black curve) growth trajectories for these 25 girls. It can be seen that, the fitted trajectories fit the observed data very well.

7 Discussion

In this paper we have proposed an estimation procedure for nonparametrically estimating the unknown gradient function of a first order autonomous differential equation over a finite domain, when the trajectories are strictly monotone. In this section, we discuss the asymptotic rate optimality of the proposed estimator. We show that, if the estimators of the gradient function g are restricted to a class of uniformly Lipschitz function, the optimal rate for estimation of g , i.e., of the order $n^{-2p/(2p+3)}$, is the same as the optimal rate for estimation of the derivative of $X \equiv X_g$ based on model (3) in terms of the L^2 loss. We conjecture that the Lipschitz requirement on the estimator of g is not necessary and the minimax rate for estimation of g is indeed of the order $n^{-2p/(2p+3)}$.

In order to make this statement precise, we first specify the function class for g as

$$\mathcal{G} = \{g : D \rightarrow \mathbb{R}_+ : c_0 \leq g \leq c_1; |g'| \leq c_2; g \in C^p(D)\} \quad (47)$$

where $0 < c_0 < c_1 < \infty$ and $0 < c_2 < \infty$ are constants. Define the class of uniformly Lipschitz functions

$$\mathcal{L} = \{h : D \rightarrow \mathbb{R} : |h(x) - h(y)| \leq c_4|x - y| \text{ for all } x, y \in D\}$$

where $c_4 \in (0, \infty)$ depends on (at least as large as) c_2 in (47). If $g \in \mathcal{G}$, then we have $X_g \in C^{p+1}([0, 1])$ and $X'_g \in C^p([0, 1])$. In addition, we assume the observation model (3) with the noise $\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma_\varepsilon^2)$.

Let δ be as in Section 2. By the condition $c_0 \leq g \leq c_1$, we know that there exist $0 < c_0(\delta) < c_1(\delta) < \infty$ such that $c_0(\delta) \leq X_g(t) \leq c_1(\delta)$ for all $t \in [\delta, 1 - \delta]$, for all $g \in \mathcal{G}$. Define, $\|f\|_{2,\delta} = (\int_\delta^{1-\delta} (f(t))^2 dt)^{1/2}$. Then there are constants $c_2(\delta), c_3(\delta) > 0$ such that for any given estimator $\hat{g} \in \mathcal{L}$ of g ,

$$\begin{aligned} c_2(\delta) \|\hat{g} \circ X_g - g \circ X_g\|_{2,\delta}^2 &\leq \int_{X_g(\delta)}^{X_g(1-\delta)} |\hat{g}(u) - g(u)|^2 du \\ &\leq c_3(\delta) \|\hat{g} \circ X_g - g \circ X_g\|_{2,\delta}^2. \end{aligned} \quad (48)$$

Observe that $g \circ X_g = X'_g$.

On the other hand, since $X_g \in C^{p+1}([0, 1])$, there exists an estimator \hat{X}_{op} with the property that, given $\epsilon > 0$, there exists constant $K_1(\epsilon) > 0$ such that

$$\sup_{g \in \mathcal{G}} \mathbb{P}(\|\hat{X}_{op} - X_g\|_{2,\delta}^2 > K_1(\epsilon) n^{2(p+1)/(2p+3)}) < \epsilon \quad (49)$$

for all $n \geq N_1(\epsilon)$.

We define the estimator $\tilde{X}' := \hat{g} \circ \hat{X}_{op}$ for X'_g . Then, by triangle inequality,

$$\begin{aligned} \|\hat{g} \circ X_g - g \circ X_g\|_{2,\delta} &= \|\hat{g} \circ X_g - X'_g\|_{2,\delta} \\ &\geq \|\tilde{X}' - X'_g\|_{2,\delta} - \|\hat{g} \circ \hat{X}_{op} - \hat{g} \circ X_g\|_{2,\delta} \\ &\geq \|\tilde{X}' - X'_g\|_{2,\delta} - c_4 \|\hat{X}_{op} - X_g\|_{2,\delta}, \end{aligned} \quad (50)$$

where, in the last step we have used the fact that $\widehat{g} \in \mathcal{L}$.

Since $X'_g \in C^p([0, 1])$, the minimax rate of estimation of X'_g in terms of the L^2 loss $\|\cdot\|_{2,\delta}^2$ is of the order $n^{-2p/(2p+3)}$. This can be derived directly for g restricted to \mathcal{G} by only slightly modifying the arguments in Stone (1982). Combining this fact with (48), (49) and (50), we obtain that there exists $K_2 > 0$, such that

$$\liminf_{n \rightarrow \infty} \inf_{\widehat{g} \in \mathcal{L}} \sup_{g \in \mathcal{G}} \mathbb{P} \left(\int_{X_g(\delta)}^{X_g(1-\delta)} |\widehat{g}(u) - g(u)|^2 du > K_2 n^{-2p/(2p+3)} \right) > 0.$$

In other words, as long as \widehat{g} is uniformly Lipschitz, the rate $n^{-2p/(2p+3)}$ is a lower bound on the rate for estimating g in terms of the L^2 -loss. We note that, the requirement $\widehat{g} \in \mathcal{L}$ can be relaxed by only requiring that this holds with probability approaching one as $n \rightarrow \infty$. The latter is satisfied by the estimator we proposed. Thus, combining with Theorem 3.2, we deduce that the optimal rate of estimation of g is $n^{-2p/(2p+3)}$ for $p \geq 4$.

8 Appendix

In this section, we provide technical details for the proofs of the main results. Specifically, in Appendix A, we present results on perturbation analysis of differential equations that are central to controlling the bias in the estimates. In Appendix B, we verify that condition (vi) of **A2** is satisfied by a B-spline basis of sufficiently high order. In Appendix C, we prove Lemma 4.2. Further technical details are given in the Supplementary Material.

Appendix A : Properties of sample trajectories and their derivatives

Throughout this subsection, with slight abuse of notation, we use $X(\cdot)$ to mean $X(\cdot; \beta)$, unless otherwise noted.

Since $X(\cdot)$ satisfies the ODE

$$X(t) = x_0 + \int_0^t \sum_{k=1}^M \beta_k \phi_k(X(s)) ds, \quad t \in [0, 1], \quad (\text{A.1})$$

differentiating with respect to β we obtain the the linear differential equations:

$$\frac{d}{dt} X^{\beta_r}(t) = X^{\beta_r}(t) \sum_{k=1}^M \beta_k \phi'_k(X(t)) + \phi_r(X(t)), \quad X^{\beta_r}(0) = 0, \quad (\text{A.2})$$

for $r = 1, \dots, M$, where $X^{\beta_r}(t) := \frac{\partial X(t)}{\partial \beta_r}$. The Hessian of $X(\cdot)$ with respect to β is given by the matrix $(X^{\beta_r, \beta_{r'}})_{r, r'=1}^M$, where $X^{\beta_r, \beta_{r'}}(t) := \frac{\partial^2 X(t)}{\partial \beta_r \partial \beta_{r'}}$, which satisfies the system of ODEs, for $r, r' = 1, \dots, M$:

$$\begin{aligned} & \frac{d}{dt} X^{\beta_r, \beta_{r'}}(t) \\ &= \left[X^{\beta_r, \beta_{r'}}(t) \sum_{k=1}^M \beta_k \phi'_k(X(t)) + X^{\beta_r}(t) \phi'_{r'}(X(t)) \right. \\ & \quad \left. + X^{\beta_{r'}}(t) \phi'_r(X(t)) + X^{\beta_r}(t) X^{\beta_{r'}}(t) \sum_{k=1}^M \beta_k \phi''_k(X(t)) \right], \quad X^{\beta_r, \beta_{r'}}(0) = 0. \end{aligned} \quad (\text{A.3})$$

With $a := X(\delta)$ and $X^a(t)$ denoting $\frac{\partial}{\partial a} X(t)$, we also have

$$\frac{d}{dt} X^a(t) = g'_\beta(X(t)) X^a(t), \quad X^a(\delta) = 1. \quad (\text{A.4})$$

Note that (A.2), (A.3) and (A.4) are linear differential equations. If the function $g_\beta := \sum_{k=1}^M \beta_k \phi_k$ is positive on the domain then the gradients of the trajectories can be solved explicitly as follows.

$$X^{\beta_r}(t) = g_\beta(X(t)) \int_{x_0}^{X(t)} \frac{\phi_r(x)}{(g_\beta(x))^2} dx. \quad (\text{A.5})$$

$$\begin{aligned} X^{\beta_r, \beta_{r'}}(t) &= g_\beta(X(t)) \int_0^t \frac{1}{g_\beta(X(s))} \left[X^{\beta_r}(s) \phi'_{r'}(X(s)) + \phi'_r(X(s)) X^{\beta_{r'}}(s) \right] ds \\ &+ g_\beta(X(t)) \int_0^t \frac{1}{g_\beta(X(s))} X^{\beta_r}(s) X^{\beta_{r'}}(s) g''_\beta(X(s)) ds. \end{aligned} \quad (\text{A.6})$$

and

$$X^a(t) = \frac{g_\beta(X(t))}{g_\beta(a)}, \quad t \in [\delta, 1 - \delta]. \quad (\text{A.7})$$

Now we summarize approximations of various relevant quantities. The following result on the perturbation of the solution path in an initial value problem due to a perturbation in the gradient function is derived from Deuffhard and Bornemann (2002).

Proposition A.1. *Consider the initial value problem:*

$$x' = f(t, x), \quad x(t_0) = x_0, \quad (\text{A.8})$$

where $x \in \mathbb{R}^d$. On the augmented phase space Ω , say, let the mappings f and δf be continuous and continuously differentiable with respect to the state variable. Assume that for $(t_0, x_0) \in \Omega$, the initial value problem (A.8), and the perturbed problem

$$x' = f(t, x) + \delta f(t, x), \quad x(t_0) = x_0,$$

have the solutions x and $\bar{x} = x + \delta x$, respectively. If f is such that $\|f_x(t, \cdot)\|_\infty \leq \chi(t)$ for a function $\chi(\cdot)$ bounded on $[t_0, t_1]$, and $\|\delta f(t, \cdot)\|_\infty \leq \tau(t)$ for some nonnegative function $\tau(\cdot)$ on $[t_0, t_1]$, then

$$\|\delta x(t)\| \leq \int_{t_0}^t \exp\left(\int_s^t \chi(u) du\right) \tau(s) ds, \quad \text{for all } t \in [t_0, t_1].$$

We use the above result to compute bounds for the trajectories and their derivatives corresponding to the different values of the parameter β in a neighborhood of the point β^* . In order to keep the exposition simple, we assume that $g_\beta(x) = g_\beta(x_{1,M})$ for $x > x_{1,M}$ and $g_\beta(x) = g_\beta(x_{0,M})$ for $x < x_{0,M}$ with a differentiability requirement at the points $x_{0,M}$ and $x_{1,M}$.

Our aim is to show that the range of the trajectories $X(t; \beta, x_{0,\delta})$ is contained in the set $D_0 = [x_{0,M}, x_{1,M}]$, for all $t \in [\delta, 1 - \delta]$ and for all $\beta \in \mathcal{B}(\alpha_n) := \{\beta : \|\beta - \beta^*\| \leq \alpha_n\}$. Let $\gamma_n = \max\{\sup_{x \in D} |g_{\beta^*}(x) - g(x)|, \sup_{x \in D_0} |g_{\beta^*}(x) - g_\beta(x)|\}$. Then $\gamma_n = O(M^{-p}) + O(\alpha_n M^{1/2})$. Also, let $\xi_n = \max_{j=0,1} |\hat{x}_j - x_{j,\delta}|$. As in the proof of Proposition 3.1, we can easily show that $[x_{0,\delta}, x_{1,\delta}] \subset [x_{0,M}, x_{1,M}]$ for sufficiently large

M . On the other hand, by using the perturbation bound given by Proposition A.1 progressively over small subintervals of the interval $[\delta, 1 - \delta]$, it can be shown that

$$\sup_{\beta \in \mathcal{B}(\alpha_n)} \sup_{t \in [\delta, 1 - \delta]} |X(t; \beta, x_{0,\delta}) - X_g(t; x_{0,\delta})| \leq C_1 \gamma_n + C_2 \xi_n,$$

for appropriate positive constants C_1, C_2 that depend on the value of g and g' on the interval $[x_0, x_1]$. Now, using Lemma 3.1, the condition on α_n as given in Theorem 3.1, and the definitions of \hat{x}_j , $x_{j,\delta}$ and $x_{j,M}$, for $j = 0, 1$, we conclude that for large enough M , the range of $X(t; \beta, x_{0,\delta})$ is contained in D_0 for all $t \in [\delta, 1 - \delta]$ and for all $\beta \in \mathcal{B}(\alpha_n)$. The scenario is depicted in Figure A.1, where the dashed curves indicate the envelop of the trajectories $X(t; \beta, x_{0,\delta})$, while the solid curve indicates the trajectory $X_g(t; x_{0,\delta})$.

Next, we provide bounds for trajectories and their derivatives. In the following, $\|\cdot\|_\infty$ is used to denote the sup-norm over $D_0 = [x_{0,M}, x_{1,M}]$. First, by **A2** we have the following:

$$\|g_\beta^{(j)} - g_{\beta^*}^{(j)}\|_\infty = O(\|\beta - \beta^*\| M^{j+1/2}) \quad j = 0, 1, 2, \quad (\text{A.9})$$

where $g^{(j)}$ and $g_{\beta^*}^{(j)}$ denote the j -th derivative of g and g_{β^*} , respectively. Next, again from **A2**, for M large enough, solutions $\{X(t; \beta) : t \in [\delta, 1 - \delta]\}$ exist for all β such that $\|\beta - \beta^*\| \leq \alpha_n$. This also implies that the solutions $X^{\beta_r}(\cdot; \beta)$ and $X^{\beta_r, \beta_{r'}}(\cdot; \beta)$ exist on $[\delta, 1 - \delta]$ for all β such that $\|\beta - \beta^*\| \leq \alpha_n$, since they follow linear differential equations where the coefficient functions depend on $X(t; \beta)$. Moreover, by *Gronwall's lemma* (Deuflhard and Bornemann, 2002), (A.9) and the fact that $\|g_{\beta^*}^{(j)}\|_\infty = O(1)$ for $j = 0, 1, 2$ (again by **A2**).

Hence, if $\|\beta - \beta^*\| M^{3/2} = o(1)$, then using Proposition A.1, the fact that $\|g_{\beta^*}^{(j)}\|_\infty = O(1)$ for $j = 0, 1, 2$, and the expressions for the ODEs for the partial derivatives, we obtain (almost surely):

$$\|X(\cdot; \beta) - X_g(\cdot)\|_\infty = O(M^{-p}). \quad (\text{A.10})$$

The same technique can be used to prove the following:

$$\|X(\cdot; \beta) - X(\cdot; \beta^*)\|_\infty = O(\|\beta - \beta^*\| M^{1/2}) \quad (\text{A.11})$$

$$\max_{1 \leq r \leq M} \|X^{\beta_r}(\cdot; \beta) - X^{\beta_r}(\cdot; \beta^*)\|_\infty = O(M^{-1/2}) \quad (\text{A.12})$$

$$\max_{1 \leq r \leq M} \|X^{\beta_r}(\cdot; \beta) - X^{\beta_r}(\cdot; \beta^*)\|_\infty = O(\|\beta - \beta^*\| M) \quad (\text{A.13})$$

$$\max_{1 \leq r, r' \leq M} \|X^{\beta_r, \beta_{r'}}(\cdot; \beta) - X^{\beta_r, \beta_{r'}}(\cdot; \beta^*)\|_\infty = O(1) \quad (\text{A.14})$$

$$\max_{1 \leq r, r' \leq M} \|X^{\beta_r, \beta_{r'}}(\cdot; \beta) - X^{\beta_r, \beta_{r'}}(\cdot; \beta^*)\|_\infty = O(\|\beta - \beta^*\| M^{3/2}) \quad (\text{A.15})$$

whenever $\|\beta - \beta^*\| M^{3/2} = o(1)$.

To illustrate the key arguments, we prove (A.12) and (A.13). First, (A.12) follows by (A.5), and the fact that $\|\phi_r\|_\infty = O(M^{1/2})$ and is supported on an interval of length $O(M^{-1})$. In fact it holds for all β such that $\|\beta - \beta^*\| M^{3/2} = o(1)$. Next, note that the function ϕ_r is Lipschitz with Lipschitz constant $O(M^{3/2})$ and is supported on an interval of length $O(M^{-1})$. Since (A.2) is a linear differential equation, using Proposition A.1 with $\delta f(t, x)$ given by

$$x \left[g'_\beta(X(t; \theta, \beta)) - g'_{\beta^*}(X(t; \beta^*)) \right] + \phi_r(X(t; \beta)) - \phi_r(X(t; \beta^*))$$

we obtain (A.13) by using (A.11) and the following facts: $\sup_{t \in [\delta, 1 - \delta]} |X^{\beta_r}(t; \beta)| = O(M^{-1/2})$ for all $\beta \in \Omega(\alpha_n)$; $\|g''_\beta\|_\infty = O(\alpha_n M^{5/2})$; $\|g'_\beta - g'_{\beta^*}\|_\infty = O(\alpha_n M^{3/2})$; and $\alpha_n M^{3/2} = o(1)$.

Appendix B : Verification of (vi) of A2 for B-spline basis

In this subsection, we verify that the condition (vi) of **A2** is satisfied if $\{\phi_{k,M}\}_{k=1}^M$ is a normalized B-spline basis with equally spaced knots on $[x_{0,M}, x_{1,M}]$ and of order $d \geq \max\{3, p-1\}$. In particular, we show that the rate of approximation of $X(t)$ by $X(t; \beta^*)$ with a carefully chosen $\beta = \beta^*$ satisfies the requirement that $\sup_{t \in [\delta, 1-\delta]} |X(t) - X(t; \beta^*)| = O(M^{-(p+1)})$ and the conditions $\sup_{x \in [x_{0,M}, x_{1,M}]} |g^{(j)}(x) - g_{\beta^*}^{(j)}(x)| = O(M^{-p+j})$ for $j = 0, 1, 2$. The result is proved through the following lemmas proved in SM.

Lemma A.1. *Suppose that $\{\phi_{k,M}\}_{k=1}^M$ has combined support $[x_{0,\delta}, x_{1,\delta}] = [X(\delta), X(1-\delta)]$ and satisfies (ii)–(v) of **A2** and β^* furthermore has the property that*

$$\sup_{x \in [x_{0,\delta}, x_{1,\delta}]} \left| \int_{x_{0,\delta}}^x \frac{g(u) - g_{\beta^*}(u)}{g(u)} du \right| = a_M \quad (\text{A.16})$$

such that $c_0 M^{-(p+1)} \leq a_M \ll M^{-p-\epsilon}$, uniformly in M , for some $\epsilon \in (0, 1]$ and some $c_0 > 0$. Then, if $X(\delta; \beta^*) = X(\delta)$, there exists $C > 0$ such that

$$\sup_{t \in [\delta, 1-\delta]} |X(t) - X(t; \beta^*)| \leq C a_M. \quad (\text{A.17})$$

Lemma A.2. *Suppose that **A1** holds with $p \geq 2$. Let $\{\phi_{k,M}\}_{k=1}^M$ denotes the normalized B-spline basis of order $\geq (p-1)$ with equally spaced knots on the interval $[x_{0,M}, x_{1,M}]$. Then there exists a $\beta^* \in \mathbb{R}^M$ such that $g_{\beta^*} = \sum_{k=1}^M \beta_k^* \phi_{k,M}$ satisfies*

$$\sup_{x \in [x_{0,\delta}, x_{1,\delta}]} \left| \int_{x_{0,\delta}}^x \frac{g(u) - g_{\beta^*}(u)}{g(u)} du \right| = O(M^{-(p+1)}). \quad (\text{A.18})$$

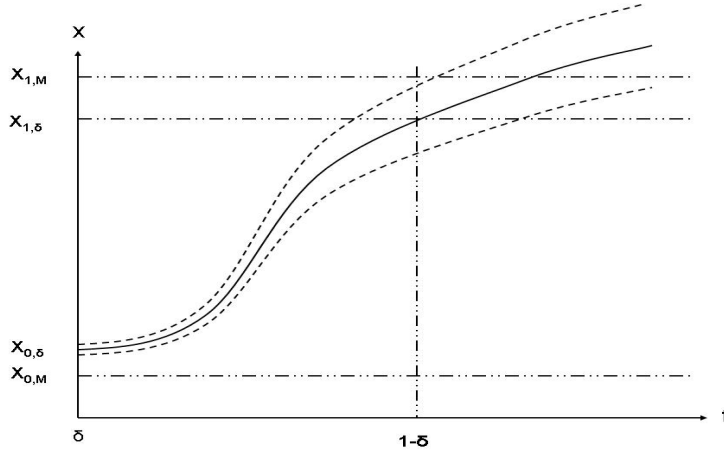


Figure A.1: Schematic diagram of the trajectory $X_g(t; x_{0,\delta})$ (solid curve) and the envelop of trajectories $X(t; \beta, x_{0,\delta})$ (boundaries indicated by dashed curves).

Appendix C : Proof of Lemma 4.2

By a Taylor expansion we have, for $j = 1, \dots, n$,

$$\begin{aligned} X(T_j; \beta) - X(T_j; \beta^*) &= X^\beta(T_j; \beta^*)^T (\beta - \beta^*) \\ &\quad + (X^\beta(T_j; \tilde{\beta}(T_j)) - X^\beta(T_j; \beta^*))^T (\beta - \beta^*), \end{aligned}$$

where $\|\tilde{\beta}(T_j) - \beta^*\| \leq \|\beta - \beta^*\|$ for all j . From this, it follows that, for all $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$,

$$\begin{aligned} \Gamma_n(\beta, \beta^*) &\geq \frac{3}{4}(\beta - \beta^*)^T \left[\frac{1}{n} \sum_{j=1}^n X^\beta(T_j; \beta^*) X^\beta(T_j; \beta^*)^T \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right] (\beta - \beta^*) \\ &\quad - 3 \|\beta - \beta^*\|^2 \frac{1}{n} \sum_{j=1}^n \|X^\beta(T_j; \tilde{\beta}(T_j)) - X^\beta(T_j; \beta^*)\|^2 \mathbf{1}_{[\delta, 1-\delta]}(T_j), \quad (\text{A.19}) \end{aligned}$$

where we have used $|2ab| \leq a^2/4 + 4b^2$. Using Proposition 3.1 and Lemma A.3 (stated below) we conclude, given $\eta > 0$, there exists $C_{10}(\eta) > 0$ such that,

$$\begin{aligned} &(\beta - \beta^*)^T \left[\frac{1}{n} \sum_{j=1}^n X^\beta(T_j; \beta^*) X^\beta(T_j; \beta^*)^T \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right] (\beta - \beta^*) \\ &\geq C_{10}(\eta) \frac{1}{M^2} \|\beta - \beta^*\|^2 \end{aligned}$$

for all $\beta \in \mathcal{A}_M(\alpha_n, \bar{\alpha}_n)$, with probability at least $1 - n^{-\eta}$. Now, another application of the Mean Value Theorem yields that for $T_j \in [\delta, 1 - \delta]$,

$$\begin{aligned} &\|X^\beta(T_j; \tilde{\beta}(T_j)) - X^\beta(T_j; \beta^*)\|^2 \\ &\leq \|\tilde{\beta}(T_j) - \beta^*\|^2 \|X^{\beta\beta^T}(T_j; \beta^*)\|_F^2 \\ &\quad + \|\tilde{\beta}(T_j) - \beta^*\|^2 \sum_{1 \leq k, k' \leq M} |X^{\beta_k, \beta_{k'}}(T_j; \tilde{\beta}^k(T_j)) - X^{\beta_k, \beta_{k'}}(T_j; \beta^*)|^2, \end{aligned}$$

where $\|\cdot\|_F$ denotes the Frobenius norm, and $\|\tilde{\beta}^k(T_j) - \beta^*\| \leq \|\tilde{\beta}(T_j) - \beta^*\|$ for all $1 \leq k \leq M$ and $1 \leq j \leq n$. Now, using (A.14) and (A.15), and combining the last three displays, we get (35).

Lemma A.3. Suppose that **A1–A4** hold. Let

$$\bar{G}_{*n} := \frac{1}{F_T(1 - \delta) - F_T(\delta)} \frac{1}{n} \sum_{j=1}^n X^\beta(T_j; \beta^*) (X^\beta(T_j; \beta^*))^T \mathbf{1}_{[\delta, 1-\delta]}(T_j).$$

Then, given $\eta > 0$, there exists constants $c'_1(\eta), c'_2(\eta) > 0$ such that, with probability $1 - n^{-\eta}$, uniformly in $\gamma \in \mathbb{S}^{M-1}$,

$$\gamma^T \bar{G}_{*n} \gamma \geq \gamma^T G_* \gamma - c'_1(\eta) \sqrt{\gamma^T G_* \gamma} \sqrt{\frac{M \log n}{n}} \geq c'_2(\eta) M^{-2}. \quad (\text{A.20})$$

Proof of Lemma A.3

Let $\mathbf{v}_j = X^\beta(T_j; \beta^*)$. Define $D(\gamma) = \gamma^T(\bar{G}_{*n} - G_*)\gamma$. Notice that

$$\frac{1}{(F_T(1-\delta) - F_T(\delta))} \mathbb{E}_T[\mathbf{v}_j \mathbf{v}_j^T \mathbf{1}_{[\delta, 1-\delta]}(T_j)] = \mathbb{E}_{\tilde{T}}[\mathbf{v}_j \mathbf{v}_j^T] = G_*,$$

where the first expectation is with respect to the distribution of T_1 and the second with respect to that of \tilde{T}_1 . Hence, we can write $D(\gamma) = n^{-1} \sum_{j=1}^n u_j(\gamma)$ where

$$u_j(\gamma) = \gamma^T \left(\mathbf{v}_j \mathbf{v}_j^T \frac{\mathbf{1}_{[\delta, 1-\delta]}(T_j)}{F_T(1-\delta) - F_T(\delta)} - \mathbb{E}_{\tilde{T}}[\mathbf{v}_j \mathbf{v}_j^T] \right) \gamma.$$

Note that, the random variables $u_j(\gamma)$ have zero conditional mean, are uniformly bounded, and are independent. Moreover, the functions $u_j(\gamma)$ are differentiable functions of γ . Then, since by (A.12), $u_j(\gamma)$'s are uniformly bounded by some $K_1 > 0$,

$$\text{Var} \left(\sum_{j=1}^n u_j(\gamma) \right) = \sum_{j=1}^n \mathbb{E}[(u_j(\gamma))^2] \leq K_1 \sum_{j=1}^n \mathbb{E}|u_j(\gamma)| \leq 2K_1 n \gamma^T G_* \gamma.$$

Thus, by Bernstein's inequality, for every $v > 0$ and $\gamma \in \mathbb{S}^{M-1}$,

$$\mathbb{P} \left(\left| \sum_{j=1}^n u_j(\gamma) \right| > v \right) \leq 2 \exp \left(- \frac{v^2/2}{2K_1 n \gamma^T G_* \gamma + K_1 v/3} \right).$$

On the other hand, by (12), $\gamma^T G_* \gamma \geq cM^{-2}$ for some $c > 0$. By this, and the condition that $M^3 = o(n/\log n)$, it is easy to see that $\sqrt{\gamma^T G_* \gamma} \gg \sqrt{M \log n/n}$. Thus, using an entropy argument as in the proof of (42), we conclude that given $\eta > 0$ there exists $c'_1(\eta) > 0$ such that

$$\mathbb{P} \left(\sup_{\gamma \in \mathbb{S}^{M-1}} \frac{|n^{-1} \sum_{j=1}^n u_j(\gamma)|}{\sqrt{\gamma^T G_* \gamma}} \leq c'_1(\eta) \sqrt{\frac{M \log n}{n}} \right) > 1 - n^{-\eta}. \quad (\text{A.21})$$

Recalling the definition of $D(\gamma)$, and again using the fact that $\gamma^T G_* \gamma \geq cM^{-2}$ and $M^3 = o(n/\log n)$, (A.20) follows from (A.21).

References

- [1] Brunel, N. J.-B. (2008). Parameter estimation of ODE's via nonparametric estimators. *Electronic Journal of Statistics* **2**, 1242-1267.
- [2] Brunk, H. D. (1970). Estimation of isotonic regression. In *Nonparametric Techniques in Statistical Inference*, Ed. Puri, M. L.
- [3] Cao, J., Fussmann, G. F. and Ramsay, J. O. (2008). Estimating a predator-prey dynamical model with the parameter cascades method. *Biometrics* **64**, 959-967.
- [4] Cao, J. and Zhao, H. (2008). Estimating dynamic models for gene regulation networks. *Bioinformatics*, **24**, 1619-1624.

- [5] Cavalier, L. (2008). Nonparametric statistical inverse problems. *Inverse Problems*, **24**, 034004.
- [6] Cavalier, L., Golubev, G. K., Lepskii, O. and Tsybakov, A. B. (2004). Block thresholding and sharp adaptive estimation in severely ill-posed inverse problems. *Theory Probability Application*, **48**, 426–446.
- [7] Chen, J. and Wu, H. (2008a). Estimation of time-varying parameters in deterministic dynamic models with application to HIV infections. *Statistica Sinica* **18**, 987–1006.
- [8] Chen, J. and Wu, H. (2008b). Efficient local estimation for time-varying coefficients in deterministic dynamic models with applications to HIV-1 dynamics. *Journal of American Statistical Association* **103**, 369–384.
- [9] Deuffhard, P. and Bornemann, F. (2002). *Scientific Computing with Ordinary Differential Equations*. Springer.
- [10] Donoho, D. L. (1995). Nonlinear solution of linear inverse problems by wavelet-vaguelette decomposition. *Applied and Computational Harmonic Analysis*, **2**, 102–126.
- [11] Erickson, R. O. (1976). Modelling of plant growth. *Annual Review of Plant Physiology*, **27**, 407–434.
- [12] Fan, J. and Gijbels, I. (1996). *Local Polynomial Modeling and Its Applications*. Chapman & Hall.
- [13] Gardner, T. S., di Bernardo, D., Lorenz, D. and Collins, J. J. (2003). Inferring genetic networks and identifying compound mode of action via expression profiling. *Science*, **301**, 102–105.
- [14] Gasser, T. and Müller, H.-G. (1984). Estimating regression functions and their derivatives by the kernel method. *Scandinavian Journal of Statistics*, **11**(3), 171–185.
- [15] Hauspie, R. C., Wachholder, A., Baron, G., Cantraine, F., Susanne, C. and Graffar, M. (1980). A comparative study of the fit of four different functions to longitudinal data of growth in height of Belgian girls. *Annals of Human Biology*, **7**(4), 347–358.
- [16] Hall, P. and Ma, Y. (2014). Quick and easy one-step parameter estimation in differential equations. *Journal of the Royal Statistical Society, Series B*. To appear.
- [17] Jank, W. and Shmueli, G. (2006). Functional data analysis in electronic commerce research. *Statistical Science*, **61**, 155–166.
- [18] Johnstone, I. M., Kerkycharian, G., Picard, D. and Raimondo, M. (2004). Wavelet deconvolution in a periodic setting. *Journal of Royal Statistical Society, Series B*, **66**, 1–27.
- [19] Kelly, C. and Rice, J. (1990). Monotone smoothing with application to doseresponse curves and the assessment of Synergism. *Biometrics*, **46**, 1071–1085.
- [20] Liu, B. and Müller, H.-G. (2009). Estimating derivatives for samples of sparsely observed functions, with application to online auction dynamics. *Journal of the American Statistical Association*, **104**, 704–716.
- [21] Ljung, L. and Glad, T. (1994). *Modeling of Dynamical Systems*. Prentice Hall.
- [22] Mammen, E. (1991). Estimating a smooth monotone regression function. *Annals of Statistics*, **19**, 724–740.

- [23] Miao, H., Dykes, C., Demeter, L. M. and Wu, H. (2009). Differential equation modeling of HIV viral fitness experiments : model identification, model selection, and multimodel inference. *Biometrics*, **65**, 292–300.
- [24] Milani, S. (2000). Kinetic models for normal and impaired growth. *Annals of Human Biology*, **27**(1), 1–18.
- [25] Mitrinovic, D. S., Pecaric, J. E. and Fink, A. M. (1991). *Inequalities Involving Functions and Their Integrals and Derivatives*. Kluwer Academic Publishers.
- [26] Müller, H.-G., Stadtmüller, U. and Schmitt, T. (1987). Bandwidth choice and confidence intervals for derivatives of noisy data. *Biometrika*, **74**, 743–749.
- [27] Nicol, F. (2013). Functional Principal Component Analysis of Aircraft Trajectories. *ISIATM 2013, 2nd International Conference on Interdisciplinary Science for Innovative Air Traffic Management, Toulouse, France*.
- [28] Olhede, S. (2007). Discussion on the paper by Ramsay, Hooker, Campbell and Cao. *Journal of Royal Statistical Society, Series B*, **69**, 772–779.
- [29] Paul, D., Peng, J. and Burman, P. (2011). Semiparametric modeling of autonomous nonlinear dynamical systems with application to plant growth. *Annals of Applied Statistics*, **5**, 2078–2108.
- [30] Perthame, B. (2007). *Transport Equations in Biology*. Birkhäuser.
- [31] Poyton, A. A., Varziri, M. S., McAuley, K. B., McLellan, P. J. and Ramsay, J. O. (2006). Parameter estimation in continuous dynamic models using principal differential analysis. *Computers & Chemical Engineering* **30**, 698–708.
- [32] Qi, X. and Zhao, H. (2010). Asymptotic efficiency and finite-sample properties of the generalized profiling estimation of parameters in ordinary differential equations. *Annals of Statistics*, **38**, 435–481.
- [33] Ramsay, J. O. (1988). Monotone regression splines in action (with discussions). *Statistical Science*, **3**, 425–461.
- [34] Ramsay, J. O. (1998). Estimating smooth monotone functions. *Journal of the Royal Statistical Society, Series B*, **60**, 365–375.
- [35] Ramsay, J. O., Hooker, G., Campbell, D. and Cao, J. (2007). Parameter estimation for differential equations: a generalized smoothing approach. *Journal of the Royal Statistical Society, Series B* **69**, 741–796.
- [36] Ramsay, J. O. and Silverman, B. W. (2005). *Functional Data Analysis*. Springer, New York.
- [37] Sacks, M. M., Silk, W. K. and Burman, P. (1997). Effect of water stress on cortical cell division rates within the apical meristem of primary roots of maize. *Plant Physiology* **114**, 519–527.
- [38] Silk, W. K., and Erickson, R. O. (1979). Kinametics of plant growth. *Journal of Theoretical Biology*, **76**, 481–501.
- [39] Stone, C. J. (1982). Optimal global rates of convergence for nonparametric regression. *Annals of Statistics*, **10**, 1040–1053.

- [40] Strogatz, S. H. (2001). *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry and Engineering*. Perseus Books Group.
- [41] Turchin, P. (2003) : *Complex Population Dynamics: A Theoretical/Empirical Synthesis*. Princeton University Press.
- [42] Tuddenham, R. D. and Snyder, M. M. (1954). Physical Growth of California Boys and Girls from Birth to Eighteen years. *University of California Publications in Child Development*, **1**, 183-364.
- [43] Varah, J. M. (1982). A spline least squares method for numerical parameter estimation in differential equations. *SIAM Journal on Scientific Computing* **3**, 28–46.
- [44] Wang, S., Jank, W., Shmueli, G. and Smith, P. (2008). Modeling price dynamics in eBay auctions using differential equations. *Journal of the American Statistical Association*, **103**, 1100–1118.
- [45] Wright, I. W. and Wegman, E. J. (1980). Isotonic, convex and related splines. *Annals of Statistics*, **8**, 1023–1035.
- [46] Wu, H., Ding, A. and DeGruttola, V. (1998). Estimation of HIV dynamic parameters. *Statistics in Medicine* **17**, 2463–2485.
- [47] Wu, H. and Ding, A. (1999). Population HIV-1 dynamics in vivo : applicable models and inferential tools for virological data from AIDS clinical trials. *Biometrics* **55**, 410–418.
- [48] Wu, H., Lu, T., Xue, H. and Liang, H. (2014). Sparse additive ordinary differential equations for dynamic gene regulatory network modeling. *Journal of the American Statistical Association*, **109**, 700–716.
- [49] Xia, X. (2003). Estimation of HIV/AIDS parameters. *Automata*, **39**, 1983–1988.
- [50] Xue, H., Miao, H. and Wu, H. (2010). Sieve estimation of constant and time-varying coefficients in nonlinear ordinary differential equation models by considering both numerical error and measurement error. *Annals of Statistics*, **38**, 2351–2387.

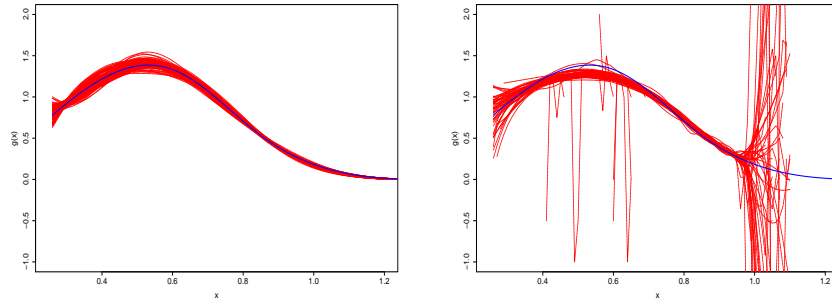


Figure 1: Simulation: Estimated gradient functions (red curves) overlaid on the true gradient function (blue curve). Left panel: proposed estimator; Right panel: two-stage estimator.

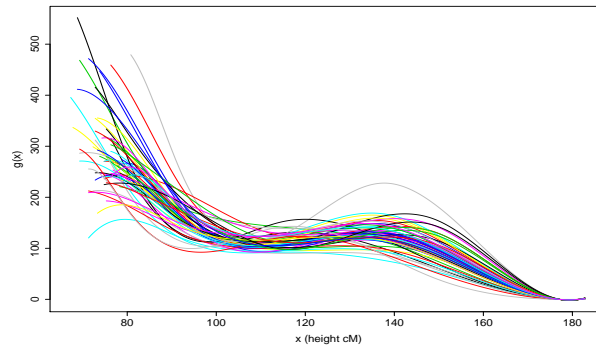


Figure 2: Berkeley Growth Data: fitted gradient functions for 54 female subjects.

Supplementary Material : “Nonparametric estimation of dynamics of monotone trajectories”

S1 Proof of Lemmas A.1 and A.2

Proof of Lemma A.1

First, we write (since $X(\delta; \beta^*) = X(\delta) = x_{0,\delta}$), for $t \in [\delta, 1 - \delta]$,

$$\begin{aligned} X(t) - X(t; \beta^*) &= \int_{\delta}^t (g(X(s)) - g_{\beta^*}(X(s; \beta^*))) ds \\ &= \int_{\delta}^t (g(X(s)) - g_{\beta^*}(X(s))) ds + \int_{\delta}^t (g_{\beta^*}(X(s)) - g_{\beta^*}(X(s; \beta^*))) ds \\ &= \int_{X(\delta)}^{X(t)} (g(u) - g_{\beta^*}(u)) \frac{du}{g(u)} + \int_{\delta}^t (g_{\beta^*}(X(s)) - g_{\beta^*}(X(s; \beta^*))) ds. \end{aligned} \quad (\text{S.1})$$

In the last step we have used $X'(s) = g(X(s))$. Since $\sup_{x \in [x_{0,\delta}, x_{1,\delta}]} |g(x) - g_{\beta^*}(x)| = O(M^{-p})$, from (A.10) we already have $\sup_{t \in [\delta, 1-\delta]} |X(t) - X(t; \beta^*)| \leq C_1 M^{-p}$ for some $C_1 > 0$. Also, $\sup_{x \in [x_{0,\delta}, x_{1,\delta}] \cap \Pi^c} \|g'_{\beta^*}(x)\| \leq C_2$ for some $C_2 > 0$, where Π denotes the set of knots for the B-spline functions. This implies, by Mean Value Theorem, that g_{β^*} is a Lipschitz function with Lipschitz constant bounded by C_2 . Thus, for all $t \in [\delta, 1 - \delta]$,

$$\begin{aligned} \int_{\delta}^t (g_{\beta^*}(X(s)) - g_{\beta^*}(X(s; \beta^*))) ds &\leq C_2 \int_{\delta}^t |X(s) - X(s; \beta^*)| ds \\ &\leq C_1 C_2 M^{-p} (t - \delta). \end{aligned} \quad (\text{S.2})$$

Combining (A.16) and (S.2), we have

$$|X(t) - X(t; \beta^*)| \leq a_M + C_1 C_2 M^{-p} (t - \delta), \quad t \in [\delta, 1 - \delta].$$

Substituting this again in the last line of (S.2), from (S.1), we obtain

$$\begin{aligned} |X(t) - X(t; \beta^*)| &\leq a_M + C_2 \int_{\delta}^t (a_M + C_1 C_2 M^{-p} (s - \delta)) ds \\ &= a_M + C_2 a_M (t - \delta) + C_1 C_2^2 M^{-p} \frac{(t - \delta)^2}{2!} \end{aligned}$$

for $t \in [\delta, 1 - \delta]$. By induction, it follows that for all $J \geq 1$,

$$|X(t) - X(t; \beta^*)| \leq a_M \sum_{j=0}^J C_2^j \frac{(t - \delta)^j}{j!} + C_1 C_2^{J+1} \frac{(t - \delta)^{J+1}}{(J+1)!}, \quad t \in [\delta, 1 - \delta].$$

Since $c_0 M^{-(p+1)} < a_M \ll M^{-p-\epsilon}$, we obtain (A.17) by choosing J sufficiently large and recalling the expansion of $e^{C_2(t-\delta)}$, whereby we can take $C = 2e^{C_2}$.

Proof of Lemma A.2

Define $G(x) = \int_{x_{0,M}}^x g(u)du$ for $x \in [x_{0,M}, x_{1,M}]$. It is well known (cf. de Boor, 1978, ch. XII) that, for every $d \geq p$, there exists a spline $S_d(x)$ of order $d \geq p$ with equally spaced knots with spacing $O(M^{-1})$ on the interval $[x_{0,M}, x_{1,M}]$ such that

$$\sup_{x \in [x_{0,M}, x_{1,M}]} |G^{(j)}(x) - S_d^{(j)}(x)| = O(M^{-(p+1)+j}), \quad \text{for } j = 0, 1, 2. \quad (\text{S.3})$$

Now, $S_d^{(1)}(x)$ is a spline of order $d-1$ on the same set of knots and hence can be expressed as $\sum_{k=1}^M \beta_k^* \phi_{k,M}(x)$ for all $x \in [x_{0,M}, x_{1,M}]$ if $\{\phi_{k,M}\}_{k=1}^M$ is the normalized spline basis of order $d-1$ on the same set of knots. Without loss of generality, we assume that $x_{0,M} < x_{0,\delta} < x_{1,\delta} < x_{1,M}$. Then, by integration by parts, we have

$$\begin{aligned} \int_{x_{0,\delta}}^x \frac{g(u) - S_d^{(1)}(u)}{g(u)} dx &= \frac{1}{g(x)} (G(x) - S_d(x)) - \frac{1}{g(x_{0,\delta})} (G(x_{0,\delta}) - S_d(x_{0,\delta})) \\ &\quad + \int_{x_{0,\delta}}^x \frac{g'(u)}{(g(u))^2} (G(u) - S_d(u)) du. \end{aligned}$$

Since $(g(u))^{-1}$ and $g'(u)$ are bounded on D , (A.18). As a by-product, we also have from (S.3) that $\sup_{x \in [x_{0,M}, x_{1,M}]} |g^{(j)}(x) - g_{\beta^*}^{(j)}(x)| = O(M^{-p+j})$ for $j = 0, 1$. Moreover, it can be checked that if $p \geq 2$ then $\sup_{x \in [x_{0,M}, x_{1,M}] \cap \Pi_M^c} |S_d^{(3)}(x)|$ is bounded, where Π_M constitute the knot sequence, which are the potential points of non-smoothness for S_d .

S2 Proof of Theorem 3.2

Since $\hat{\beta}$ is a local minimizer of $\tilde{L}_\delta(\beta)$, it satisfies $(\partial/\partial\beta)\tilde{L}_\delta(\hat{\beta}) = 0$. Thus, applying the Mean Value Theorem coordinatewise, we have, for $k = 1, \dots, M$,

$$-\frac{\partial}{\partial\beta_k} \tilde{L}_\delta(\beta^*) = \frac{\partial^2}{\partial\beta_k \partial\beta^T} \tilde{L}_\delta(\tilde{\beta}_k)(\hat{\beta} - \beta^*) \quad (\text{S.4})$$

for some $\tilde{\beta}_k$ such that $\|\tilde{\beta}_k - \beta^*\| \leq \|\hat{\beta} - \beta^*\|$. We write

$$\begin{aligned} -\frac{\partial}{\partial\beta} \tilde{L}_\delta(\beta^*) &= 2 \sum_{j=1}^n \varepsilon_j \frac{\partial}{\partial\beta} X(T_j; \beta^*, \hat{x}_0) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ &\quad + 2 \sum_{j=1}^n (X_g(T_j) - X(T_j; \beta^*, \hat{x}_0)) \frac{\partial}{\partial\beta} X(T_j; \beta^*, \hat{x}_0) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ &=: \tilde{U}_1 + \tilde{U}_2. \end{aligned}$$

On the other hand,

$$\begin{aligned}
\frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \tilde{L}_\delta(\boldsymbol{\beta}) &= 2 \sum_{j=1}^n \frac{\partial}{\partial \boldsymbol{\beta}} X(T_j; \boldsymbol{\beta}, \hat{x}_0) \left(\frac{\partial}{\partial \boldsymbol{\beta}} X(T_j; \boldsymbol{\beta}, \hat{x}_0) \right)^T \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad - 2 \sum_{j=1}^n \varepsilon_j \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} X(T_j; \boldsymbol{\beta}, \hat{x}_0) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad - 2 \sum_{j=1}^n (X_g(T_j) - X(T_j; \boldsymbol{\beta}, \hat{x}_0)) \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} X(T_j; \boldsymbol{\beta}, \hat{x}_0) \\
&=: \tilde{S}_1(\boldsymbol{\beta}) + \tilde{S}_2(\boldsymbol{\beta}) + \tilde{S}_3(\boldsymbol{\beta}).
\end{aligned}$$

Then, we can express (S.4) in vectorial form as

$$\begin{aligned}
-\frac{\partial}{\partial \boldsymbol{\beta}} \tilde{L}_\delta(\boldsymbol{\beta}^*) &= \tilde{S}_1(\boldsymbol{\beta}^*)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) + \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_1(\tilde{\boldsymbol{\beta}}_k) - \tilde{S}_1(\boldsymbol{\beta}^*))(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) \\
&\quad + \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_2(\tilde{\boldsymbol{\beta}}_k) + \tilde{S}_3(\tilde{\boldsymbol{\beta}}_k))(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*),
\end{aligned}$$

where \mathbf{e}_k denotes the vector in \mathbb{R}^M with 1 in k -th coordinate and zero elsewhere. From this, we get the expansion

$$\begin{aligned}
\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* &= \left(\tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \tilde{U}_1 + \left(\tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \tilde{U}_2 \\
&\quad - \left(\tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_1(\tilde{\boldsymbol{\beta}}_k) - \tilde{S}_1(\boldsymbol{\beta}^*))(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) \\
&\quad - \left(\tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_2(\tilde{\boldsymbol{\beta}}_k) + \tilde{S}_3(\tilde{\boldsymbol{\beta}}_k))(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*). \tag{S.5}
\end{aligned}$$

Let $S_l(\boldsymbol{\beta})$ be the counterpart of $\tilde{S}_l(\boldsymbol{\beta})$, $l = 1, 2, 3$ once we replace the initial condition \hat{x}_0 by $x_{0,\delta}$. It is then easily verified that for some $C_{11} > 0$,

$$\begin{aligned}
\left\| \frac{1}{n} \tilde{S}_1(\boldsymbol{\beta}^*) - \frac{1}{n} S_1(\boldsymbol{\beta}^*) \right\| &\leq C_{11} M^{1/2} |\hat{x}_0 - x_{0,\delta}| (1 + M^{1/2} |\hat{x}_0 - x_{0,\delta}|) \\
&= O_P(M^{1/2} (\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)}). \tag{S.6}
\end{aligned}$$

This implies in particular, by Proposition 3.1 and Lemma A.3 that

$$\left\| \left(\frac{1}{n} \tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \right\| \leq C_{12} M^2 + O_P(M^{1/2} (\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)}) = O(M^2) \tag{S.7}$$

for M satisfying the condition (13). Thus

$$\text{Var} \left(\left(\tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1} \tilde{U}_1 \mid \mathbf{T}, \hat{x}_0 \right) = \frac{\sigma_\varepsilon^2}{n} \left(\frac{1}{n} \tilde{S}_1(\boldsymbol{\beta}^*) \right)^{-1},$$

and hence,

$$\text{Trace} \left[\text{Var} \left(\left(\tilde{S}_1(\beta^*) \right)^{-1} \tilde{U}_1 \mid \mathbf{T}, \hat{x}_0 \right) \right] = O_P \left(\frac{\sigma_\varepsilon^2 M^3}{n} \right) \quad (\text{S.8})$$

by (S.7). On the other hand, from

$$\sup_{t \in [\delta, 1-\delta]} |X(t; \beta^*, \hat{x}_0) - X(t; \beta^*)| = O(|\hat{x}_0 - x_{0,\delta}|) = O_P((\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)})$$

and $\sup_{t \in [\delta, 1-\delta]} |X_g(t) - X(t; \beta^*)| = O(M^{-(p+1)})$, and the form of \tilde{U}_2 , we have, with an application of Proposition S.1 (stated below) that

$$\begin{aligned} \left\| \left(\tilde{S}_1(\beta^*) \right)^{-1} \tilde{U}_2 \right\| &\leq 2 \left\| \left(\frac{1}{n} \tilde{S}_1(\beta^*) \right)^{-1} \right\|^{1/2} \\ &\quad \cdot \sup_{t \in [\delta, 1-\delta]} (|X(t; \beta^*, \hat{x}_0) - X(t; \beta^*)| + |X_g(t) - X(t; \beta^*)|) \\ &= O_P(M^{-p}) + O_P(M(\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)}). \end{aligned} \quad (\text{S.9})$$

Now, using arguments analogous to those used in the proof of Theorem 3.1, and the bounds (A.12)–(A.15), we can show that the maximum of the norms of the matrices $\left(\tilde{S}_1(\beta^*) \right)^{-1} \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_1(\tilde{\beta}_k) - \tilde{S}_1(\beta^*))$ and $\left(\tilde{S}_1(\beta^*) \right)^{-1} \sum_{k=1}^M \mathbf{e}_k \mathbf{e}_k^T (\tilde{S}_l(\tilde{\beta}_k))$, for $l = 2, 3$, is $O_P(M^3 \alpha_n)$, which is $o_P(1)$ by the condition on M . From this, and (S.8) and (S.9), the result (16) follows.

Proposition S.1. *Suppose that B be an $p \times n$ matrix such that BB^T is invertible. Let y be an $n \times 1$ vector. Then $\| (BB^T)^{-1} B y \| \leq (\| (BB^T)^{-1} \|)^{1/2} \| y \|$.*

Proposition S.1 follows immediately by using singular value decomposition of B .

S3 Rate of convergence of the two-stage estimator

First, define

$$\begin{aligned} \mathbf{W} &= \frac{1}{n} \sum_{j=1}^n \phi(X(T_j)) \phi(X(T_j))^T \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\ \hat{\mathbf{W}} &= \frac{1}{n} \sum_{j=1}^n \phi(\hat{X}(T_j)) \phi(\hat{X}(T_j))^T \mathbf{1}_{[\delta, 1-\delta]}(T_j). \end{aligned}$$

Then, using the fact that $X'(t) = g(X(t))$ and $g_{\beta^*}(x) = \phi(x)^T \beta^*$, we have

$$\begin{aligned}
\tilde{\beta} &= \hat{\mathbf{W}}^{-1} \frac{1}{n} \sum_{j=1}^n \phi(X(T_j)) \phi(X(T_j))^T \beta^* \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad + \hat{\mathbf{W}}^{-1} \frac{1}{n} \sum_{j=1}^n (g(X(T_j)) - g_{\beta^*}(X(T_j))) \phi(X(T_j)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad + \hat{\mathbf{W}}^{-1} \frac{1}{n} \sum_{j=1}^n (\hat{X}'(T_j) - X(T_j)) \phi(X(T_j)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&\quad + \hat{\mathbf{W}}^{-1} \frac{1}{n} \sum_{j=1}^n \hat{X}'(T_j) (\phi(\hat{X}(T_j)) - \phi(X(T_j))) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \\
&= \beta^* + \mathbf{W}^{-1} (\mathbf{W} - \hat{\mathbf{W}}) \hat{\mathbf{W}}^{-1} \mathbf{W} \beta^* + R_1 + R_2 + R_3,
\end{aligned} \tag{S.10}$$

where R_1 , R_2 and R_3 are the expressions in the second, third and fourth lines after the first equality.

We check that the following bounds hold with probability tending to 1 for any given sequence $M \rightarrow \infty$ as $n \rightarrow \infty$.

$$\max\{\|\mathbf{W}\|, \|\mathbf{W}^{-1}\|\} = O(1). \tag{S.11}$$

$$\|\hat{\mathbf{W}} - \mathbf{W}\| = O(M^2 (\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)} \sqrt{\log n}). \tag{S.12}$$

$$\|\mathbf{W}^{-1} \left(\frac{1}{n} \sum_{j=1}^n (g(X(T_j)) - g_{\beta^*}(X(T_j))) \phi(X(T_j)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right)\| = O(M^{-p}). \tag{S.13}$$

$$\begin{aligned}
&\left\| \mathbf{W}^{-1} \left(\frac{1}{n} \sum_{j=1}^n (\hat{X}'(T_j) - X'(T_j)) \phi(X(T_j)) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right) \right\| \\
&= O((\sigma_\varepsilon^2/n)^{p/(2p+3)} \sqrt{\log n}).
\end{aligned} \tag{S.14}$$

$$\begin{aligned}
&\left\| \frac{1}{n} \sum_{j=1}^n \hat{X}'(T_j) (\phi(\hat{X}(T_j)) - \phi(X(T_j))) \mathbf{1}_{[\delta, 1-\delta]}(T_j) \right\| \\
&= O(M^{3/2} (\sigma_\varepsilon^2/n)^{(p+1)/(2p+3)} \sqrt{\log n}).
\end{aligned} \tag{S.15}$$

Combining these with (S.10) we obtain Proposition 3.2.

Proof of (S.11)

First, write \mathbf{W} as $\bar{\mathbf{W}} + \Delta_W$, where

$$\bar{\mathbf{W}} = \int_{\delta}^{1-\delta} \phi(X(t)) \phi(X(t))^T f_T(t) dt = \int_{X(\delta)}^{X(1-\delta)} \phi(u) (\phi(u))^T \frac{f_T(X^{-1}(u))}{g(u)} du.$$

Notice that for any $\mathbf{y} \in \mathbb{S}^{M-1}$, $\mathbf{y}^T \bar{\mathbf{W}} \mathbf{y}$ lies in the interval

$$\left(\int_{X(\delta)}^{X(1-\delta)} (\mathbf{y}^T \phi(u))^2 du \right) \left[\frac{\min_{s \in [0,1]} f_T(s)}{\max_{u \in [X(0), X(1)]} g(u)}, \frac{\max_{s \in [0,1]} f_T(s)}{\min_{u \in [X(0), X(1)]} g(u)} \right]$$

from which it follows that $\max\{\|\bar{\mathbf{W}}\|, \|\bar{\mathbf{W}}^{-1}\|\} = O(1)$ (uniformly in M) by the property of the B-spline basis (Schumaker, 2007). The result then follows from the fact (derived along the line of Lemma A.3) that $\|\Delta_W\| \leq c(\eta)M\sqrt{\log n/n} = o(1)$ with probability $1 - n^{-\eta}$ for any given $\eta > 0$.

Proof of (S.12)

This follows from the observation that

$$\begin{aligned} & \|\hat{\mathbf{W}} - \mathbf{W}\| \\ & \leq \frac{1}{n} \sum_{j=1}^n (\|\phi(\hat{X}(T_j))\| + \|\phi(X(T_j))\|) \|\phi(\hat{X}(T_j)) - \phi(X(T_j))\| \mathbf{1}_{[\delta, 1-\delta]}(T_j), \end{aligned}$$

and then using Mean Value Theorem, followed by condition (iii) of **A2**, and finally invoking (20), we get the result.

Proof of (S.13)

Here, if we denote the vector inside $\|\cdot\|$ by γ , then we have

$$\mathbf{W}\gamma = \frac{1}{n} \sum_{j=1}^n (g(X(T_j)) - g_{\beta^*}(X(T_j)))\phi(X(T_j))\mathbf{1}_{[\delta, 1-\delta]}(T_j).$$

Taking inner product with γ , applying Cauchy-Schwarz inequality on the right and then using (A.9), we have $\gamma^T \mathbf{W}\gamma \leq c_9 M^{-p} \sqrt{\gamma^T \mathbf{W}\gamma}$ for some $c_9 > 0$. Hence, by (S.11), we have the result.

Proof of (S.14)

It is similar to that of (S.13) and uses (21) rather than (A.9).

Proof of (S.15)

It uses similar arguments as in the proof of (S.12).

S4 Sub-Gaussian random variables

We summarize a few facts about sub-Gaussian random variables. The following is a restatement of Lemma 5.5 of Vershynin (2011).

Lemma S.1. *A random variable ξ is sub-Gaussian, if any of the following equivalent conditions hold.*

- (1) $\mathbb{E}(e^{\xi^2/K_1^2}) < \infty$ for some $0 < K_1 < \infty$
- (2) $(\mathbb{E}(|\xi|^q))^{1/q} \leq K_2 \sqrt{q}$ for all $q \geq 1$, for some $0 < K_2 < \infty$.

If moreover, $\mathbb{E}(\xi) = 0$, then the following is equivalent to (1) and (2).

- (3) $\mathbb{E}(e^{t\xi}) \leq e^{t^2 K_3^2}$ for all $t \in \mathbb{R}$, for some $0 < K_3 < \infty$.

Define the *sub-Gaussian norm* of a random variable ξ to be

$$\|\xi\|_{\psi_2} := \sup q^{-1/2} (\mathbb{E}|\xi|^q)^{1/q}. \quad (\text{S.16})$$

Clearly, by Lemma S.1, ξ is a sub-Gaussian random variable if and only if $\|\xi\|_{\psi_2} < \infty$.

One of the useful characteristics of sub-Gaussianity is that it is preserved under linear combinations. Specifically, we have the following result.

Lemma S.2. (*Lemma 5.9 in Vershynin (2011)*). Suppose that X_1, \dots, X_n are independent sub-Gaussian random variables and $b_1, \dots, b_n \in \mathbb{R}$ are nonrandom quantities. Then $\sum_{i=1}^n b_i X_i$ is sub-Gaussian and

$$\left\| \sum_{i=1}^n b_i X_i \right\|_{\psi_2}^2 \leq C \sum_{i=1}^n b_i^2 \|X_i\|_{\psi_2}^2 \quad (\text{S.17})$$

for some $C > 0$.

The result follows easily from the equivalent characterizations in Lemma S.1, specifically, by using the moment generating function. The following simple corollary is very useful for our applications.

Corollary S.1. Suppose that X_1, \dots, X_n are independent random variables with $\max_{1 \leq i \leq n} \|X_i\|_{\psi_2} \leq K < \infty$. Then $\sum_{i=1}^n b_i X_i$ is sub-Gaussian and

$$\left\| \sum_{i=1}^n b_i X_i \right\|_{\psi_2}^2 \leq CK^2 \left(\sum_{i=1}^n b_i^2 \right) \quad (\text{S.18})$$

for some $C > 0$.

The following (Proposition 5.10 in Vershynin (2011)) is a version of *Hoeffding's inequality* for sub-Gaussian random variables.

Lemma S.3. Let ξ_1, \dots, ξ_n be independent random variables satisfying $\mathbb{E}(\xi_i) = 0$, and let $K := \max_{1 \leq i \leq n} \|\xi_i\|_{\psi_2} < \infty$. Then for any $b_1, \dots, b_n \in \mathbb{R}$ we have

$$\mathbb{P} \left(\left| \sum_{i=1}^n b_i \xi_i \right| > t \right) \leq e \exp \left(-\frac{ct^2}{K^2 \sum_{i=1}^n b_i^2} \right), \quad \text{for all } t > 0, \quad (\text{S.19})$$

for some $c > 0$.

References

- [1] de Boor, C. (1978). *A Practical Guide to Splines*. Springer-Verlag.
- [2] Schumaker, L. (2007). *Spline Functions : Basic Theory*. Cambridge University Press.
- [3] Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv:1011.3027*.

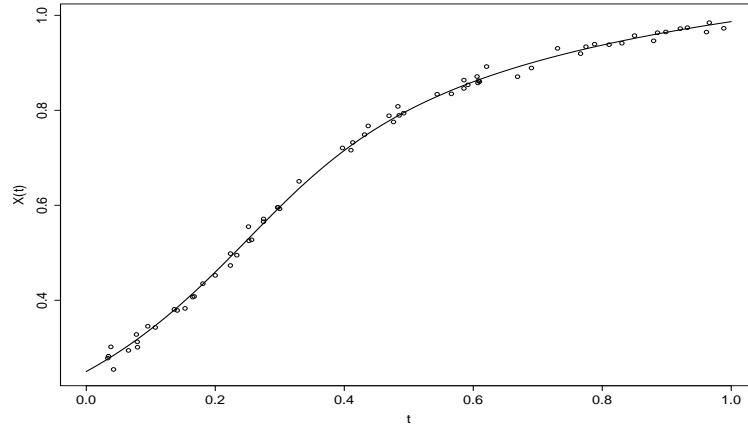


Figure S.1: Simulation: True trajectory and sample observations for one replicate.

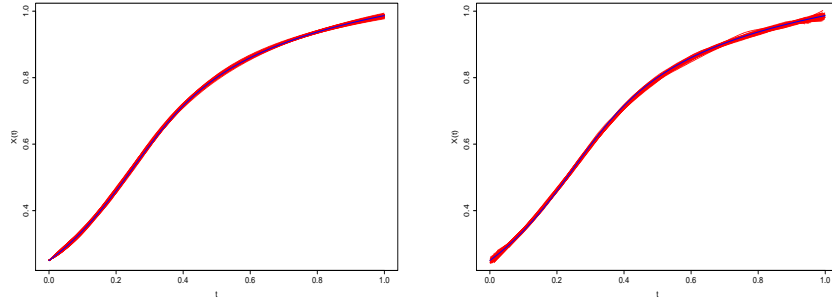


Figure S.2: Simulation: Estimated trajectory $\hat{X}(\cdot)$ (red curves) overlayed on the true trajectory $X(\cdot)$ (blue curve). Left panel: proposed estimator; Right panel: estimator from the 1st stage smoothing of the two-stage procedure.

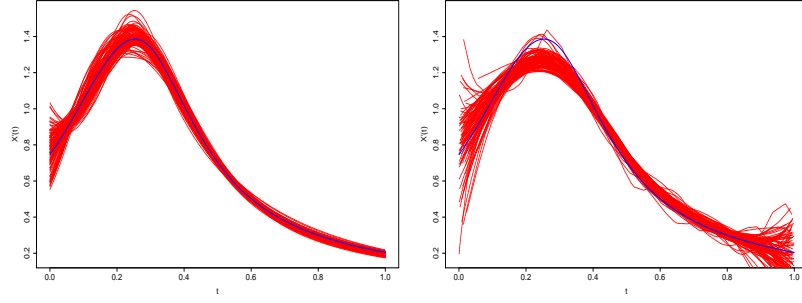


Figure S.3: Simulation: Estimated derivative of the trajectory $\hat{X}'(\cdot)$ (red curves) overlaid on the true derivative of the trajectory $X'(\cdot)$ (blue curve). Left panel: proposed estimator; Right panel: estimator from the 1st stage smoothing of the two-stage procedure.

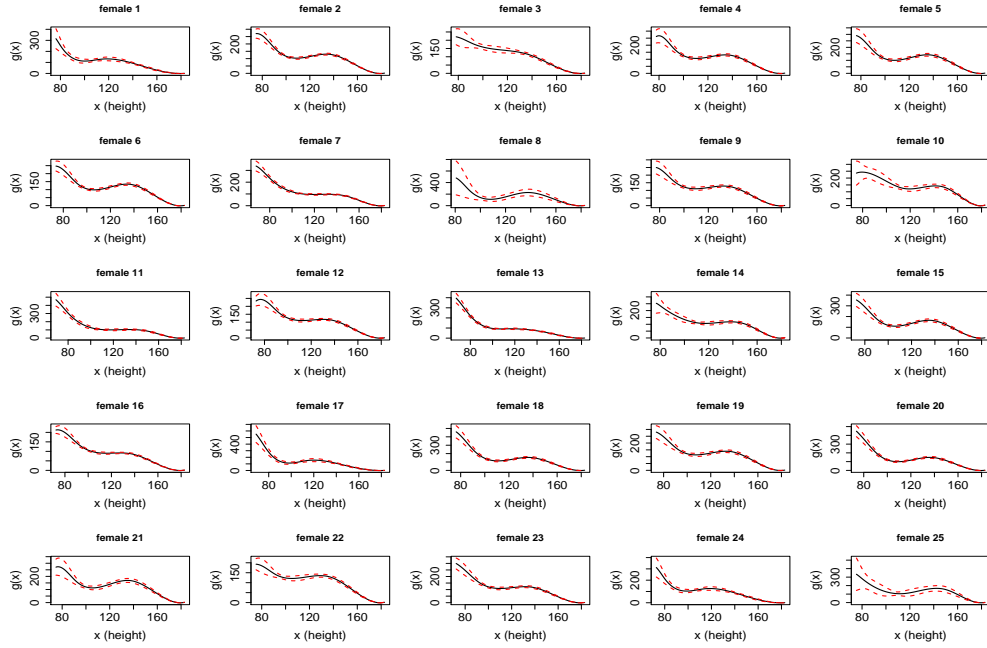


Figure S.4: Berkeley Growth Data: fitted gradient functions (black curve) for 25 female subjects with two-standard-error bands (red broken lines).

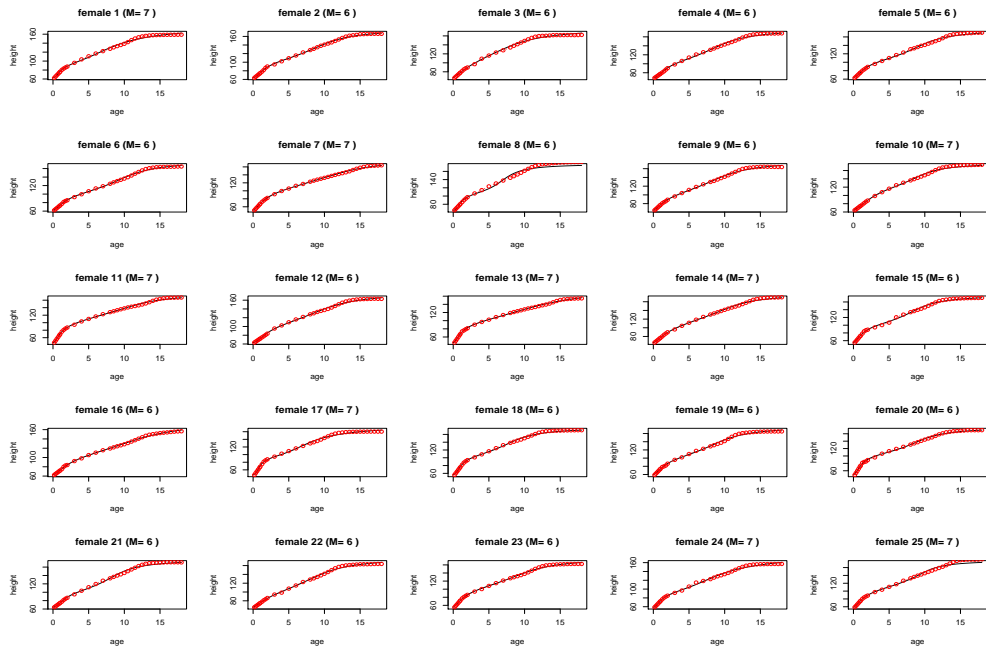


Figure S.5: Berkeley Growth Data: observed (red dots) and fitted (black curve) growth trajectories for 25 female subjects.